24

25

# Supporting Appendix 1: Additional Methods and Results for Computer-Assisted Initial Diagnosis of Rare Diseases

Rui Alves[1,*,§], Marc Piñol, [2,*], Jordi Vilaplana[2], Ivan Teixido[2], Joaquim Cruz[1], Jorge Comas[1], Ester Vilaprinyo[1], Albert Sorribas[1], Francesc Solsona[2,§],

[1]Departament d'Informàtica i Enginyeria Industrial, Universitat de Lleida, Av. Jaume II nº 69, 25001 LLeida, Spain.
[2]Departament de Ciències Mèdiques Bàsiques & IRBLleida, Universitat de Lleida, Montserrat Roig nº 2, 25008 LLeida, Spain.
e-mails:      wrrzag666@gmail.com
              jordi@diei.udl.cat
              iteixido@diei.udl.cat
              joaquimcruz92@gmail.com
              jorgecomasp@gmail.com
              evilaprinyo@cmb.udl.cat
              albert.sorribas@cmb.udl.cat
              francesc@diei.udl.cat
              ralves@cmb.udl.cat

*These two authors contribute equally for the study
§Corresponding authors:

    Francesc Solsona: Departament d'Informàtica i Enginyeria Industrial, Universitat de Lleida, Av. Jaume II nº 69,
                        25001 LLeida, Spain.
                        email: francesc@diei.udl.cat
                        tel: 00 34 973 70 27 35
                        fax: 00 34 973 70 27 02

    Rui Alves: Departament de Ciències Mèdiques Bàsiques & IRBLleida, Universitat de Lleida,
                        Edifici de Recerca Biomèdica I, Av. Rovira Roure 80, 25198 LLeida, Spain.
                        email: ralves@cmb.udl.cat
                        tel: 00 34 973 70 24 25
                        fax: 00 34 973 70 24 26

**65**     **TECHNOLOGY UNDERLYING THE RARE DISEASE DISCOVERY PROTOTYPE**

**66**     The technology behind the web application is the GRAILS framework, a web application

**67**     framework built for the Java Virtual Machine that uses the Groovy programming language. GRAILS

**68**     uses a MVC (Model View Controller) pattern that allows for a full integration between the model

**69**     (and the database) and the view (user interface). With built-in database access and modeling, it

**70**     enables easy abstraction and decoupling between these two parts of the application, permitting

**71**     easy database migrations. This also helps hiding database complexity and access to information in

**72**     an object-oriented way. JQuery and Ajax were also used in order to provide dynamic web

**73**     capabilities like autocomplete. A powerful front-end framework for faster and easier web

**74**     development (Twitter Bootstrap) was included, streamlining the styling and design of the web

**75**     interface. Built for the JVM, this framework also enables easy integration with Java packages,

**76**     plugins and wrappers.

**77**     The database design provides a welcome positive side-effect: it is trivial to keep the database up-

**78**     to-date. A periodic download of the ORPHANET data every three months, followed by upload of

**79**     that data to our database can be done in minutes, facilitating that RDD is kept up to date and

**80**     usable over the long run.  Currently, the database has a total of 6 915 diseases and 2 110

**81**     symptoms. There is a total of 101 840 records representing relations between symptoms and

**82**     diseases.

**83**     We note that the symptoms-disease association file from the ORPHANET dataset are re-curated by

**84**     us in order to ensure that the automated processing of the xml file made available by ORPHANET

**85**     is done without mistakes. Although it does not always happen, in some versions of the xml files we

**86**     downloaded had one or more tags that were not properly closed. In addition, in the earlier

**87**     versions of the files, the symptoms had not yet been fully converted to their synonymous terms in

88    the HPO (Human Phenotype Ontology)[1]. We performed a script analysis to identify those terms

89    that were not in HPO and transformed them into their HPO synonyms. In the last 3 versions of the

90    ORPHANET xml file, we found that HPO nomenclature has been fully implemented.

91

92    **CHOOSING THE APPROPRIATE PREDICTION METHODS**

93    Other classification approaches to predicting rare diseases based on symptoms were tested. First,

94    we tested an additional ranking function that takes into consideration how frequently each

95    symptom is thought to be associated with the disease. This information is provided in the

96    ORPHANET dataset, which associates qualitative frequency information to a symptom, when it is

97    associated to a disease (Very Frequent, Frequent, and Occasional). This tested function as the

98    form:

99    $$DS'_i = 1 - \frac{\sum_{j=1}^{n} \delta_j}{Max[S_{user}, S_{Disease\ i}]}$$    Eq. A1

100   In Eq. A1 $S_{user}$ represents the number of symptoms provided by the user, $S_{Disease\ i}$ represents the

101   number of symptoms of disease $i$ stored in the database, and $Max[S_{user}, S_{Disease\ i}]$ represents the

102   largest number between $S_{user}$ and $S_{Disease\ i}$. $n$ represents the number of symptoms that are different

103   between the set submitted by the user and the set associated to any given rare disease in the

104   database. $\delta_j$ measures the qualitative frequency at which symptom $j$ has been found to associate

105   to disease $i$ in the past (see above). Given that there were only three categorical frequency

106   associations (Very Frequent, Frequent, Occasional), $\delta_j$ was considered to have one of three values.

107   $\delta_j = 1$ if the symptom is either very frequently associated to the disease $i$ or is a symptom that is

108   provided by the user; $\delta_j = 0.75$ if the symptom is frequently associated to the disease $i$; finally, if

109   the symptom is only occasionally associated to disease $i$, $\delta_j = 0.5$. It can be shown that ,

110     $-1 \leq DS_i \leq DS_i' \leq 1$.    However, even though $DS_i' \neq DS_i$ , the list of diseases is ranked in the

111     same order by both scores (data not shown). Because more calculations are required to estimate

112     $DS_i'$, using this score for ranking leads to slower computation. Hence, we discarded $DS_i'$.

113     Second, we also trained and tested algorithms based on Support Vector Machines, Neural

114     Networks, Bayesian Networks, Random Trees, and Random Forests. Invariably, these algorithms

115     required extensive training and prediction time, and their best performance was always about one

116     order of magnitude lower than that of the algorithm and score described in this paper. They were

117     also orders of magnitude slower in predicting the disease and required more computational

118     resources for doing so.

119     **RETROSPECTIVE STUDY OF PREVIOUSLY DIAGNOSED RARE DISEASE PATIENTS**

120     RAMEDIS is a server that provides management services for medical doctors diagnosing, treating

121     and managing rare disease patients. Its database contains short report cards with at most 3

122     sentences about 1099 patients with confirmed rare disease clinical diagnostics.

123     The information for about 60% of these patients is public, although anonymized. From

124     these approximately six hundred patients, nearly half have metabolic rare diseases that were

125     diagnosed in screening programs at a preclinical stage. Of the remainder three hundred patients,

126     one hundred and eighty seven had a confirmed clinical diagnosis associated with a report card that

127     described at least one symptom. Examples of the procedure are given

128     We took these 187 patients and reconstructed their symptoms from the individual report

129     cards. In some cases this is easy, and report cards were very clear (for example: patient with

130     seizures or hypotonia). In other cases the symptoms were vaguely described and hard to

131     reconstruct. For example, "hearing problems" or "hearing loss" could be any of the following:

132     "Conductive deafness/hearing loss", "Central deafness/hearing loss", "Sensorineural

133    deafness/hearing loss", or "Hearing loss/hypoacusia/deafness". Another example, "Infection"

134    could be any of the following: "Immunodeficiency/increased susceptibility to infections/recurrent

135    infections", "Recurrent urinary infections", "Chronic skin infection/ulcerations/ulcers/cancrum", or

136    "Repeat respiratory infections". In these cases, we opted for including all possibilities rather than

137    eliminating the symptom. This decision was made because eliminating the symptom would have

138    meant discarding additional patients from an already small set, as all reported symptoms were

139    often ambiguous. An example of two report cards and their processing is shown in Supporting

140    Figure 1. The patients, their symptoms, and their clinically confirmed diagnosis can be manually

141    accessed and compiled from the RAMEDIS website. Supporting Figure 2 plots the accumulated

142    frequency of the score for the correct (and best) prediction.

143    **BENCHMARKING THE RARE DISEASE DISCOVERY PROTOTYPE**

144    The rare disease prediction algorithm was extensively benchmarked to evaluate the effect of

145    absent and unrelated symptoms on diagnostic precision. In addition, we also tested how the

146    changes in the ORPHANET dataset could affect the results. These benchmarks relied on several

147    sets of tests, all run using Stochastic Monte-Carlo simulations.

148    **Aggregated effects of unreported and unrelated symptoms on prediction accuracy of the Rare**

149    **Disease Discovery Algorithm**

150    The first benchmark test was done by generating several random sets of 10 000 patients, each

151    with all the symptoms associated to a specific but randomly chosen rare disease. Then, for

152    increasing percentages of the patients in a given random set either 1, 2, 3, 4, 5, 10, or 20

153    symptoms were randomly added or deleted to create noise. Then, the noisy sets of symptoms

154    were used by the RDD algorithm to predict the rare disease that generated them. The precision $p$,

155    sensitivity *s*, and *F-Score* of the RDD prediction algorithm were calculated for each set of patients.

156    The results are summarized in Figure 2 of the main text and discussed in the main manuscript.

157    **Effects of unreported symptoms on prediction accuracy of the Rare Disease Discovery Algorithm**

158    The second benchmark test was done by again generating several random sets of 10 000 patients,

159    each with all the symptoms associated to a specific but randomly chosen rare disease. Then, for

160    increasing percentages of the patients in a given random set, either 25%, 50%, or 75% of the

161    symptoms were deleted to create noise. Finally, the noisy sets of symptoms were used by the RDD

162    algorithm to predict the rare disease that generated them. These simulations represent situations

163    where not all symptoms are known to the user during diagnosis. The precision *p*, sensitivity *s*, and

164    *F-Score* of the RDD prediction algorithm were calculated for each set of patients. The results are

165    summarized in Figure 3 of the main text and discussed in the main manuscript.

166    **Estimating significance for $DS_i$ scores and testing the performance of RDD in misdiagnosing**

167    **patients that do not suffer from rare diseases**

168    It is important to estimate how large $DS_i$ must be for a user to be sure that the set of symptoms

169    being submitted to RDD (Rare Disease Discovery) are not the result of randomly associated

170    symptoms. A third benchmark of the RDD algorithm was done to estimate this $DS_i$ value. This

171    estimation was done in following way. Consider that there are 13 698 diseases and 2 528

172    symptoms in our database. The average number of symptoms associated to a disease is 42, with a

173    standard deviation of 59. To calculate the probability that a given $DS_i$ for a set of symptoms

174    produced by a user is statistically significant we generated 10 000 random vectors of symptoms.

175    The population of the 10 000 vectors had an average number of symptoms equal to 42, with a

176    standard deviation of 59. Given that these vectors were random, by plotting $f = (1 -$

177    $Accumulated\ frequency\ of\ DS_i)$ as a function of $DS_i$ (Supporting Figure 3) we are able to

178     estimate the probability that a given score is achieved simply by choosing a random combination

179     of symptoms. This experiment estimates that a score $DS_i \geq 0.5$ has a probability lower than

180     0.0001 of being obtained by choosing a random set of symptoms. If we lower the probability to

181     0.01, then $DS_i \geq 0.25$. In fact, the median $DS_i$ score for a random choice of symptoms is less than

182     0.01.

183     **Estimating significance levels for the differences between two $ds_i$ scores**

184     In the previous section we describe an experiment that allowed us to estimate that if $DS_i>0.5$, one

185     can be 99.99% sure that the score was not obtained by choosing a random set of symptoms.

186     Another issue is that of determining how significant are the differences between two $DS_i$

187     scores for the same set of symptoms. Estimating this is much more complicated because the

188     significance will depend on the number of symptoms one submits for the prediction. A final

189     benchmark experiment was done in order to provide a best scenario estimation for how

190     statistically significant the differences between two $DS_i$ scores are.

191     In this fourth and final benchmark we performed the following Monte Carlo simulation

192     experiments. For each disease we created all possible sets of k symptoms, where $k$=1, 2, 3, 4, 5,

193     10, 20, and 50 symptoms that are associated to that disease (taking care to eliminate diseases in

194     the simulation that had less than the simulated number of symptoms). Then, for each $k$, we

195     calculated $DS_i$ for the correct disease. We call this list $DS_{i\ correct}$ In parallel, for each $k$ and for each

196     set of symptoms, we calculated $DS_i$ for all diseases that were not the one from which we had

197     extracted the set of symptoms. We call this list $DS_{i\ incorrect}$.

198     Then, for each k we created a list $\Delta DS_i$, where each element of the list corresponds is

199     obtained by subtracting quantile j of $DS_{i\ incorrect}$ from quantile j of $DS_{i\ correct}$. The results are

200     presented in Supporting Table 2 and interpreted in the following way. For the same number of

201 submitted symptoms, in the context of the disease-symptoms association matrix, the differences

202 between corresponding quantiles of the $DS_{i\ correct}$ and $DS_{i\ incorrect}$ lists provide a proxy to evaluate

203 how different two $DS_i$ scores (one correct and one incorrect) must be for that difference to be

204 significant. Thus, if users submit for example one symptom and want a certainty of 99.9% that two

205 $DS_i$ scores are different, Supporting Table 2 tells us that the two scores should differ by at least

206 0.14. How can this be interpreted? For example, the difference between the score for the most

207 highly ranked disease and that for the second best guess by RDD needs to differ by at least 0.14, if

208 one want to state that the prediction is significantly (p<0.001) better than the second best guess.

209 It is important to benchmark the performance of RDD with patients that have symptom(s)

210 associated to rare diseases, without suffering from those diseases. This is a very real scenario, as

211 many of the symptoms are common between rare and non-rare diseases. A possible test would be

212 to create synthetic patients from other diseases, adding random rare disease symptoms and

213 running RDD. However, we note that RDD only allows users to choose symptoms that have been

214 previously associated to at least one rare disease. Hence, testing RDD's performance with

215 synthetic patients from non-rare diseases is formally equivalent to generating synthetic patients

216 with random associations of rare-disease symptoms. This is the same test that was done to

217 determine significance for $DS_i$ scores. In other words, only when $DS_i$ is larger than 0.5, does RDD

218 ensure that the patient has a rare disease, with a probability higher than 0.9999.


219 **Accurate predictions in the absence of statistically significant DS$_i$ scores**

220 Taken together, the four benchmark experiments described in the main manuscript show that $DS_i$

221 decreases sharply with noise; however, even if $DS_i$ is below the statistically significant level, it can

222 still be used to accurately predict the correct rare disease, although with a lower confidence (see

223 above). For example, in Supporting Figure 4 we show Box plots of the maximum $DS_i$ scores for all

224 patients in the second benchmark test. We see that when patients have 50% absent symptoms,

225   the maximum score is still almost always above 0.5, which is the 0.0001 significance level

226   determined in benchmark 3. Only when 75% of the symptoms are absent do we get maximum $DS_i$

227   scores that are equal to or lower than 0.5 for more than 50% of the patients.

228

229   **Effect of evolving datasets:  ORPHANET dataset of December 2014 vs. ORPHANET dataset of**

230   **December 2015**

231   Given that the dataset we used is annotated by humans and evolves, we wanted to have an

232   estimate of how much the changes might affect the predictive capabilities of RDD. To achieve this

233   we repeated the tests described in all the previous subsections of "BENCHMARKING THE RARE

234   DISEASE DISCOVERY ALGORITHM" for the ORPHANET dataset of 2015. What we found was that

235   the difference in F-Score of RDD between the two sets was smaller than 3% when noise was large

236   (20 noisy symptoms) and less than 0.2% when symptoms were accurate (Supporting Figure 5). We

237   also observe that the median score of the correct prediction when 25%, 50%, or 75% of symptoms

238   are absent increased by approximately 20% when we changed the 2014 dataset for the 2015

239   dataset (Supporting Figure 6). These results suggest that the human curation of the ORPHANET

240   dataset is improving over time, which also improves the quality of the results of computer assisted

241   DDX tools that use them, as is the case of RDD.

| Patient case report | Main data |
| --- | --- |
| Patient ID | 101 |
| Diagnosis | PHENYLKETONURIA (MIM 261600) |

| Gender | m |
| --- | --- |
| Age of symptoms onset | |
| Age of diagnosis | 5 Day(s) |
| Found in newborn screening | y |
| Diagnosis confirmed | y |
| Country | Germany |
| Ethnic origin | Mother: German, Father: German |
| History | Increased phe in newborn-sreening. Early start of phe-restricted dietary treatment. |

| Patient case report | Main data |
| --- | --- |
| Patient ID | 762 |
| Diagnosis | ARGININOSUCCINIC ACIDURIA (MIM 207900) |

| Gender | f |
| --- | --- |
| Age of symptoms onset | 2 Day(s) |
| Age of diagnosis | 3 Day(s) |
| Found in newborn screening | y |
| Diagnosis confirmed | n |
| Country | Germany |
| Ethnic origin | Mother: German, Father: German |
| History | The 4th day of life, the patient was hospitalized with coma and highly increased ammonia levels. In the extended newborn screening program, elevated levels of citrulline and decreased arginine-levels were found. Psychomotor development is normal. Constant hepatopathy with hepatomegaly and increased transaminases and alkaline phosphatase. |

**A**    no usable symptoms; discarded report card

**B**    symptoms: comma; hyperammonemia; hepatopathy; hepatomegaly;
vague: hepatopathy -possible symptoms include "Abnormal hepatic
enzymes/transaminases", "Hepatitis/icterus/cholestasis", "Liver/hepatic steatosis",
"Acute hepatic failure", "Chronic hepatic failure", "Hepatoblastoma", "Liver/hepatic
abscess", "Polycystic liver disease/hepatic cysts", "Intrahepatic biliary tract
atresia/obstruction", "Congenital hepatic fibrosis", "Hepatocellular liver disease/hepatic
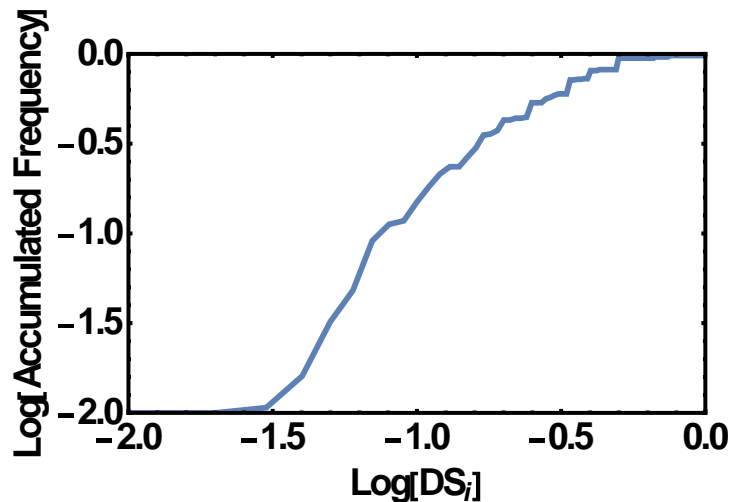failure", "Hepatitis/icterus/cholestasis"

**Supporting Figure 1** – Two examples of RAMEDIS report cards. **A**: Example of a report card that could not be used, as no symptoms were reported. **B**: Example of a report card that could be used, but had vague description of some symptoms.
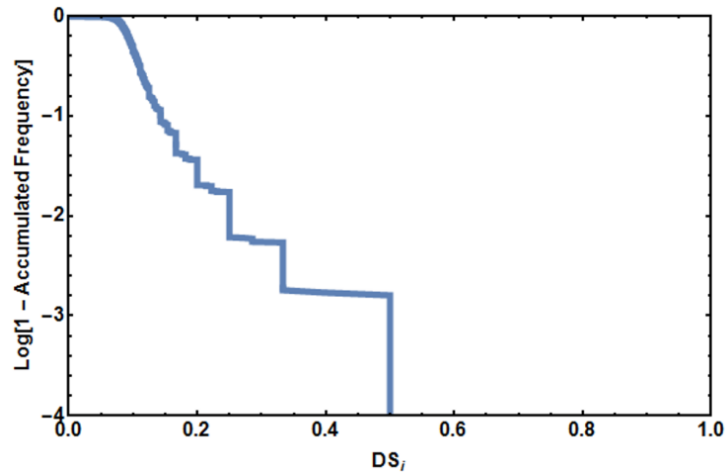


**Supporting Figure 2** – Cumulative frequency of the highest score for the retrospective study of previously diagnosed rare disease patients.
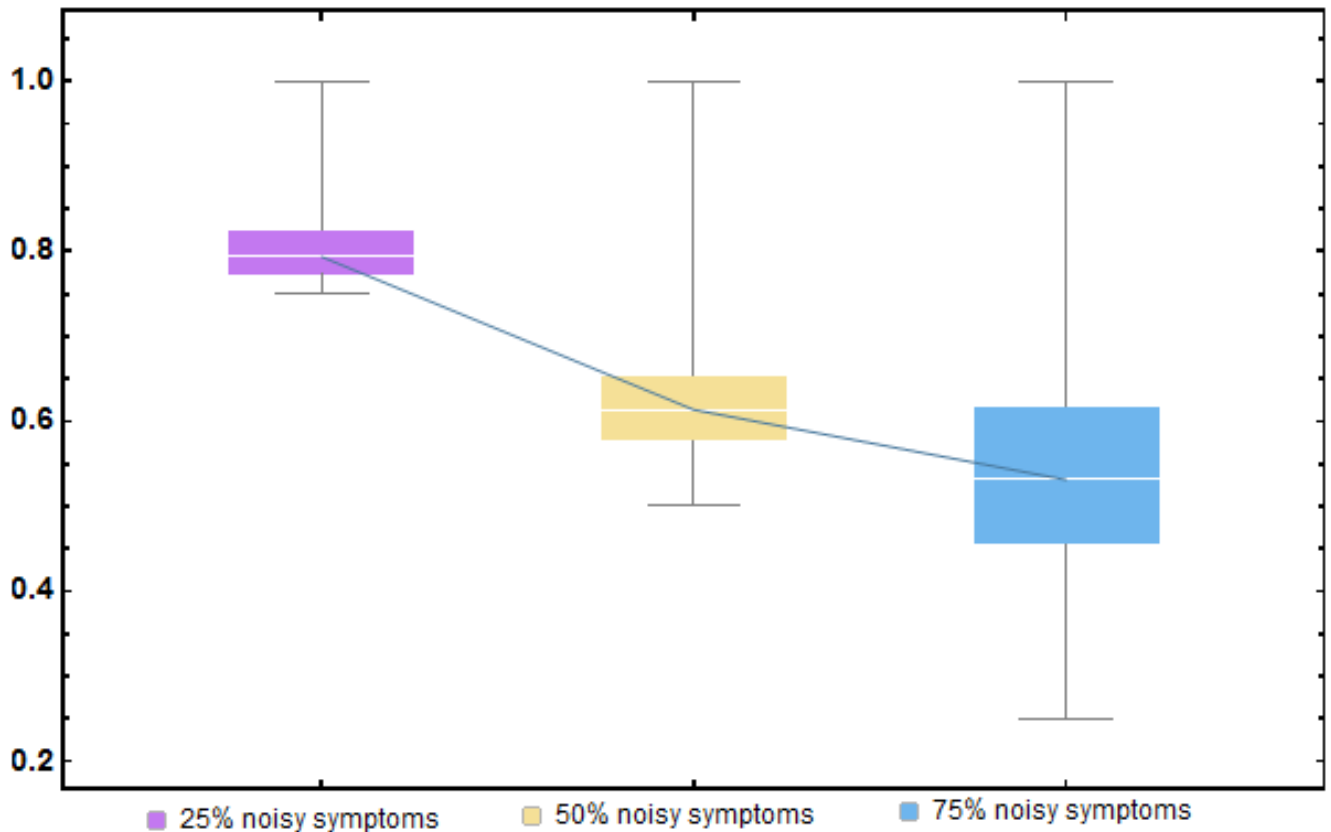
254

**Supporting Figure 3** – Estimating the probability that a score $DS_i{\leq}k$ might be obtained from a
random set of symptoms. The score $DS_i$ is represented in the x-axis, while the logarithm of 1 − the
accumulated frequency of $DS_i$ is represented in the y-axis. $DS_i$ score values higher than 0.5 have a
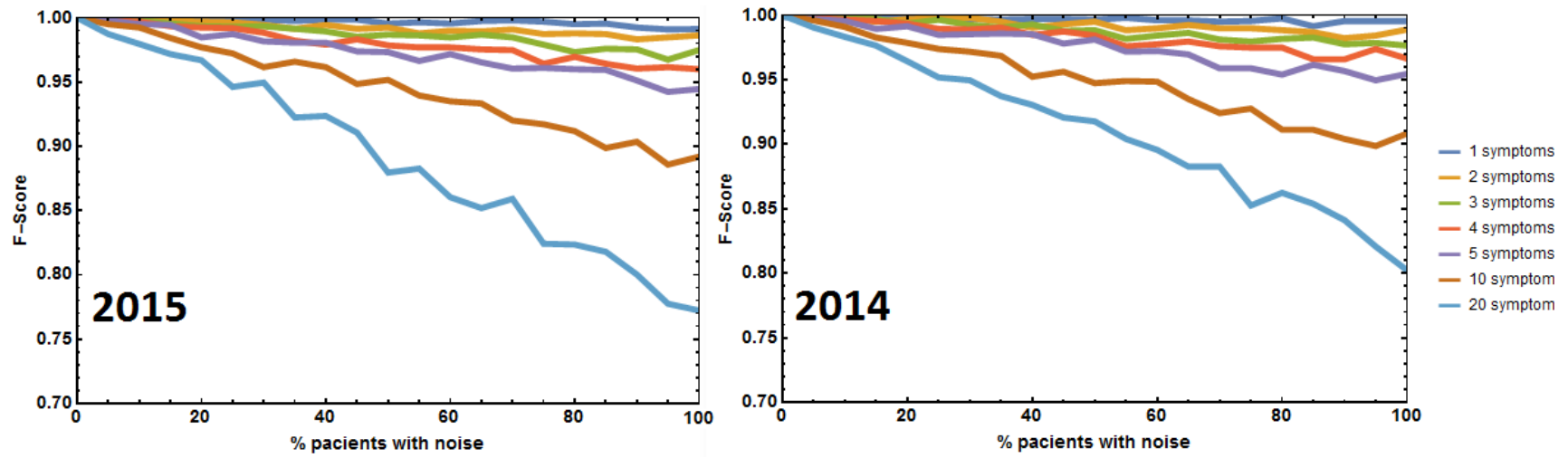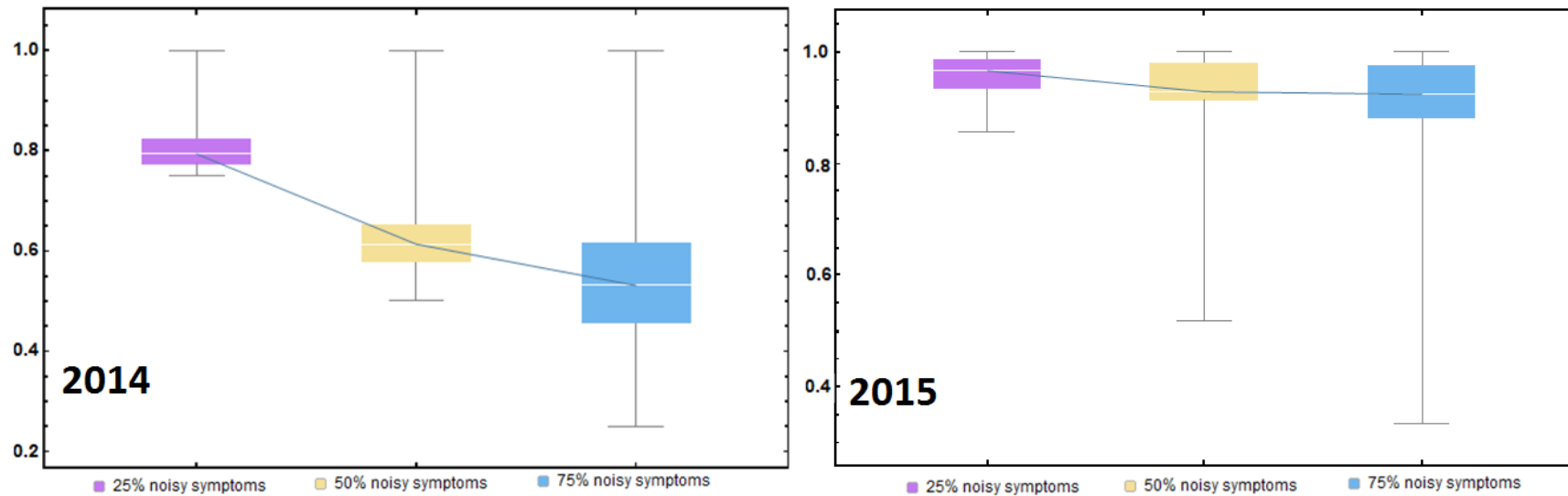probability of 0.0001 of being obtained from a random set of symptoms.

259
260



261

**Supporting Figure 4** – Effect of unreported symptoms on the maximum $DS_i$ scores for patients.
Here we present Box plots of the maximum $DS_i$ scores for all patients when 25%, 50%, or 75% of
the symptoms are absent. The median maximum scores are joined by a blue line. The boxes
indicate the 0.25 and 0.75 quartiles in each dataset.

**Supporting Figure 5** – Effect of the evolution in the ORPHANET dataset on the F1-Score of RDD. Comparison of the datasets from December 2014 and 2015. The effect of the dataset from different years is less than 3%.

**Supporting Figure 6** – Effect of the evolution in the ORPHANET dataset on the effect of unreported symptoms on the maximum $DS_i$ scores for patients. Comparison of the datasets from December 2014 and 2015. The median maximum scores are joined by a blue line. The boxes indicate the 0.25 and 0.75 quartiles in each dataset. The median scores of the newest dataset are higher than those of the 2014 dataset, indicating an improvement in the quality of the RDD predictions when symptoms are under-reported.

**Supporting Table 1 – Symptoms used to perform the experiments summarized in Tables 1 and 2.**

| Disease | Initial Symptom | Rank at First symptom | Additonal symptoms required for the appropriate disease to be ranked as 1st prediction |
|---|---|---|---|
| **Beta-Thalassemia** | Chronic skin infection/ulcerations/ulcers/cancrum | 67th | Humour troubles/anxiety/depression/apathy/euphoria/irritability<br>Anaemia |
| **Canavan Disease** | Motor deficit/trouble | 23rd | Seizures/epilepsy/absences/spasms/status epilepticus<br>Retinitis pigmentosa/retinal pigmentary changes<br>Hypotonia<br>Contractures/cramps/trismus/tetania/claudication/opisthotonos |
| **Down Syndrome** | Strabismus/squint | 244th | Sterility/hypofertility<br>Microstomia/little mouth<br>Insulin-independent/type 2 diabetes |
| **Fabry Disease** | Renal failure | 111th | Anorexia<br>Humour troubles/anxiety/depression/apathy/euphoria/irritability<br>Renal disease/nephropathy<br>Nausea/vomiting/regurgitation/merycism/hyperemesis<br>Myalgia/muscular pain<br>Thick lips<br>X-linked recessive inheritance |
| **Goldblatt Syndrome** | Hip dislocation/dysplasia/coxa valga/coxa vara/coxa plana | 81st | Delayed dentition/eruption of teeth/lack of eruption of teeth<br>Respiratory distress/dyspnea/respiratory failure/lung volume reduction |
| **Turner Syndrome** | Pigmented naevi/naevus pigmentosus/lentigo | 21st | Thin/hypoplastic toe nails |
| **Uncombable Hair Syndrome** | Albinism (hair) | 1st | Albinism (hair) |
| **Williams Syndrome** | Renal failure | 121st | Angor pectoris/myocardial infarction<br>Thin/hypoplastic toenails<br>Late puberty/hypogonadism/hypogenitalism<br>Osteosclerosis/osteopetrosis/bone condensation |
| **Yunis-Varon Syndrome** | Sternal/sternum anomalies | 7th | Cardiomyopathy/hypertrophic/dilated<br>Poorly ossified skull/calvarium<br>Absent/small toenails/anonychia of feet<br>Blepharophimosis/short palpebral fissures<br>Absent/small fingernails/anonychia of hands<br>Anteverted nares/nostrils<br>Hip dislocation/dysplasia/coxa valga/coxa vara/coxa plana<br>Hypotonia |
| **Zellweger-like Syndrome without Peroxisomal Anomalies** | High forehead | 31st | Broad nasal root<br>Expressionless face/amimia |

**Supporting Table 2** – Estimating the probability that a difference between scores $\Delta DS_i < x$ is significant at three p-value levels, when a varying number of symptoms is submitted to RDD.

| Number of symptoms | $\Delta DS_i\,(p-value < 0.01)$ | $\Delta DS_i\,(p-value < 0.005)$ | $\Delta DS_i\,(p-value < 0.001)$ |
|---|---|---|---|
| 1 | $\Delta DS_i \leq 0.01$ | $\Delta DS_i \leq 0.025$ | $\Delta DS_i \leq 0.14$ |
| 2 | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.005$ | $\Delta DS_i \leq 0.015$ |
| 3 | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.010$ |
| 4 | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.007$ |
| 5 | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.005$ |
| 10 | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.001$ |
| 20 | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.001$ |
| 50 | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.001$ | $\Delta DS_i \leq 0.001$ |