

Interpretable Deep Learning to Map Diagnostic Texts to ICD10 Codes

Aitziber Atutxa, Arantza Díaz de Ilarraza, Koldo Gojenola, Maite Oronoz, Olatz Perez-de-Viñaspre

Department of Languages and Computer Systems.

IXA Research Group: <http://ixa.eus>.

University of the Basque Country (UPV-EHU)

{aitziber.atutxa, a.diazdeilarraza, koldo.gojenola, maite.oronoz, olatz.perezdevinaspre}@ehu.eus

Abstract

Background. Automatic extraction of morbid disease or conditions contained in Death Certificates is a critical process, useful for billing, epidemiological studies and comparison across countries. The fact that these clinical documents are written in regular natural language makes the automatic coding process difficult because, often, spontaneous terms diverge strongly from standard reference terminology such as the International Classification of Diseases (ICD).

Objective. Our aim is to propose a general and multilingual approach to render Diagnostic Terms into the standard framework provided by the ICD. We have evaluated our proposal on a set of clinical texts written in French, Hungarian and Italian.

Methods. ICD-10 encoding is a multi-class classification problem with an extensive (thousands) number of classes. After considering several approaches, we tackle our objective as a sequence-to-sequence task. According to current trends, we opted to use neural networks. We tested different types

of neural architectures on three datasets in which Diagnostic Terms (DTs) have their ICD-10 codes associated.

Results and conclusions. Our results give a new state-of-the art on multilingual ICD-10 coding, outperforming several alternative approaches, and showing the feasibility of automatic ICD-10 prediction obtaining an F-measure of 0.838, 0.963 and 0.952 for French, Hungarian and Italian, respectively. Additionally, the results are interpretable, providing experts with supporting evidence when confronted with coding decisions, as the model is able to show the alignments between the original text and each output code.

Keywords: International Classification of Diseases, Electronic Health Records, Sequence-to-Sequence mapping, Neural Machine Translation

1. Introduction

Death Certificates are clinical text documents written by clinicians. They are typically accompanied by numerical codes to describe the morbid disease or conditions that led to the death of an individual. These codes originate from the International Classification of Diseases (ICD)¹. They ensure normalization upon the ways different clinicians and countries employ for writing Diagnostic Terms (DTs) to describe the same disease. For instance, the English standard description of the E141 code is *Unspecified diabetes mellitus: With ketoacidosis*, but many clinicians use different variants. Table 1 presents several variants taken from Electronic Death Certificates in different languages.

¹In this paper the 10th revision of the ICD classification will be referred to as ICD-10.

ICD10 code	French	Italian	Hungarian
E141	diabete acidocetosique decompensation diabeto-cetosique diabete acidosique diabete cetosique insuline acido-cetose diabetique acidocetose diabetique	chetoacidosi diabetic scompensat acidosi diabetic chetoacidosi diabetic acetonemia diabetic acidosi metabolic diabetic chetonemia diabetic	cukorbetegseg ketoacidosis diab versavanyodas diabetes ketoacidosis diabeteses ketoacidozis diabeteses acidosis ketoacidotikus diab

Table 1: Example of different DTs to express the E141 code in French, Italian and Hungarian.

12 The coding is useful for billing, epidemiological studies or comparison
 13 across countries. The process of assigning diagnosis codes is labor intensive,
 14 costly and error prone due to the fact that it is carried out by human coders.

15 The training of experts to encode the ICD is expensive and time-consuming.
 16 According to Lang [1], about \$25 billion per year is spent in the USA in en-
 17 coding records with ICD codes. In addition, the ICD is evolving and, there-
 18 fore, experts have to adapt to the most recent revision available. Even if
 19 standard diagnostic terminology is well defined and widely internationalized
 20 (the ICD exists in 43 languages and is used in around 117 countries), for
 21 obvious reasons, physicians usually use their own non-standard expressions
 22 that seldom agree with the standard. Hence, automatic ICD encoding is
 23 a very useful clinical application based on information extracted from text.
 24 Natural language processing methods are specially suitable to satisfactorily
 25 solve this task.

26 There have been many attempts to automatically extract ICD codes from
 27 clinical documents. This paper presents the application of different sequence-

28 to-sequence models to map DTs to ICD-10 codes [2], a large-scale multi-class
29 classification task.

30 Neural network models have revolutionized AI, especially in sequence-to-
31 sequence mapping [3, 4]. Recently, a number of open systems have emerged
32 [5, 6, 7], allowing experimentation of different approaches in a flexible way.
33 Following this paradigm, the task can be viewed as a special type of machine
34 translation (MT). In this work, we consider both input (plain text for DTs)
35 and output texts (codes) as if they were two different languages, and we try
36 to translate from one to another.

37 Our approach is general and has been applied to a varied set of languages,
38 namely, French, Italian and Hungarian, obtaining the best results in the
39 CLEF eHealth 2018 Task 1: “Multilingual Information Extraction - ICD10
40 coding”. CLEF (Conference and Labs of the Evaluation Forum) is a well-
41 known international initiative that, since 2000, has run campaigns for the
42 systematic evaluation of information access systems, playing a leading role
43 in stimulating investigation and research in a wide range of key areas in
44 the information retrieval domain. CLEF is especially interesting for the
45 comparison of approaches dealing with a specific task which is considered
46 of special interest for the community. Besides CLEF, there have been other
47 main venues related to the automatic evaluation of tasks on medical texts
48 in the last years, like the Clinical Natural Language Processing Challenges
49 (*i2b2*, *n2c2*) [8, 9] and the BioNLP-Shared Tasks [10]. The CLEF 2017
50 and 2018 eHealth shared tasks [11, 12] have been the reference forum on
51 multilingual ICD coding. The task consisted of mapping Death Certificate
52 lines containing DTs to the relevant ICD-10 codes, with eleven and fourteen

53 teams participating in 2017 and 2018, respectively. In the 2017 edition, the
54 languages under study were English and French while, in 2018, the languages
55 were Hungarian, Italian and French.

56 **2. Related work**

57 ~~The rationale behind this review is to summarize the main approaches~~
58 ~~to automatic ICD coding and their relation to the present work. Besides,~~
59 ~~we present how the scientific community deals with this problem in devoted~~
60 ~~shared tasks.~~

61 The first attempts to automatically extract ICD codes from clinical doc-
62 uments date back to the 1990s (more precisely, ICD-9). For many years
63 software for ICD-coding used Dictionary Matching and Pattern Matching
64 methods to search keywords or clusters of keywords to identify ICD codes
65 [13]. These methods were not powerful enough because the task presents a
66 complex characterization: large-scale multi-class classification, treatment of
67 non-standard language, and alignment issues between spontaneous writing
68 and ICD codes. Nowadays, there are more sophisticated approaches to tackle
69 this task, ranging from knowledge-based solutions to statistical [14, 15] and
70 deep learning ones.

71 Rule-based systems are still used with good accuracy when the terms to be
72 coded follow regular patterns, the task is limited to one DT one ICD assign-
73 ment, and the number of ICDs is quite small. Unfortunately, these conditions
74 seldom apply. An example is described in [16] where in order to classify men-
75 tal severity symptoms in psychiatric records, knowledge-based methods out-
76 perform those obtained by neural networks. In the first shared task related

77 to ICD coding, the Computational Medicine Challenge (CMC) [8], Farkas
78 and Szarvas [17] started to address this task by replacing some steps in the
79 construction of hand-crafted-systems with machine learning (ML). They re-
80 alized that manually building rules was not straightforward and it was time
81 consuming. They worked with 45 classes and not with the entire scope of the
82 ICD-10 catalog (thousands of classes). Pérez et al. [18] extracted the encod-
83 ing rules in the form of inferred Weighted Finite-State Transducers (WFST)
84 from the corpus to produce normalized alternatives to a given spontaneous
85 string and applied text similarity to select the standard string in the ICD
86 framework that best matched the normalized alternatives. In the current
87 work, a first step of normalization is not required, as neural systems learn
88 how to manage non-standard language and, in addition, instead of a one-to-
89 one (1:1) alignment, a much harder N:M term-code alignment is necessary.

90 When the range of classes to be tagged is high and the corpus is big
91 enough, machine learning based techniques have been successful [19, 20].
92 Koopman et al. [21] trained Support Vector Machine (SVM) classifiers to
93 identify cancers in a cascaded architecture, first identifying the presence of
94 a cancer and later classifying its type according to the ICD-10 classification.
95 In another work, Koopman et al. [22] used a machine learning approach
96 and keyword matching rules to identify the presence of diabetes, influenza,
97 pneumonia and HIV in Death Certificates. An SVM classifier with term-
98 based and concept-based features (i.e. SNOMED CT concepts) was trained
99 for each of the four diseases. A single classifier model was also trained for each
100 ICD-10 code representing each of these diseases. In the rule-based approach
101 a set of keywords provided by experts was used to indicate whether a Death

102 Certificate was a positive or a negative match for a particular disease. These
103 works approach the problem as document classification (Death Certificate
104 classification), instead of term encoding. In [23], separate machine learning
105 models were trained with data from unstructured text, semi-structured text
106 and structured tabular data to create a multimodal machine learning model
107 that predicts ICD-10 diagnostic codes. Recently, neural networks boosted
108 the results in many NLP tasks [24], also in the medical domain [25] including
109 ICD coding [26, 27]. Transfer learning [28] has also been used combined with
110 deep learning, as in [29] where MeSH [30] domain knowledge is transferred
111 to a deep learning model by pre-training it with MeSH datasets and fine-
112 tuning the neural network on ICD-9 datasets. The authors state that transfer
113 learning is the key element to improve ICD-9 encoding.

114 In the last CLEF eHealth ICD encoding shared tasks [11, 12], dictionaries
115 and corpora were the main resources and n-gram patterns, machine learning
116 and neural approaches the most employed methods. In the CLEF 2017 edi-
117 tion, Ebersbach et al. [31] and Zweigenbaum and Lavergne [32], along with
118 many other teams, relied on lexical resources and made use of different ML
119 methods, dictionary projections from medical ontologies, expansion of syn-
120 onyms and edit distance calculation. In this CLEF eHealth 2017 edition, only
121 Miftahutdinov and Tutubalina [33] implemented Recurrent Neural Networks
122 to assign ICD-10 codes to fragments of English Death Certificates. Their
123 system used a LSTM to map the input sequence into a vector representa-
124 tion, and then another LSTM to decode the target sequence from the vector,
125 obtaining an F-measure of 0.850.

126 In the latest CLEF 2018 edition, 12 teams out of 14 used the dictionaries

127 supplied by the organization. Cossin et al. [34] applied a dictionary-based ap-
128 proach, including a module for the detection of typos and another module for
129 synonym expansion. Regarding machine learning approaches, Almagro et al.
130 [35] implemented a supervised learning system using multilayer perceptrons
131 and a One-vs-Rest (OVR) strategy, also experimenting with IR methods. In
132 Gobeill and Ruch [36] a purely statistical instance-based approach was used
133 by indexing all training sentences to feed a k-Nearest Neighbors algorithm.
134 Neural architectures have been used to approximate the ICD encoding task.
135 Jeblee et al. [37] used an ensemble model for prediction which includes n-gram
136 matching followed by an ensemble of a convolutional neural network and a
137 recurrent neural network (RNN) encoder-decoder. Their system employed
138 embeddings learned on the data provided, as well as on language-specific
139 Wikipedia corpora.

140 Ševa et al. [38] focused on the setup and evaluation of a language-independent
141 neural architecture using a multi-language word embedding space. The au-
142 thors obtained these word representations by concatenating Italian, French
143 and Hungarian Fastext² pre-trained embeddings. Their approach builds on
144 two RNNs, modelling the extraction and classification based on Long Short-
145 Term Memories (LSTM) with an attention mechanism. Another idea to
146 initialize the networks in a recurrent neural model, in this case the weights
147 of the final network, is given in [39] where co-occurrences between ICD classes
148 in the training data and the hierarchical structure of ICD-10 are used. Note
149 that in the latter work the whole Death Certificates including the events

²<https://github.com/facebookresearch/fastText/blob/master/docs/crawl-vectors.md>

150 leading to death (diagnoses with ICD-10 codes) and the main cause of death
151 are employed. In our system co-occurrences are modeled implicitly by the
152 sequence to sequence approach.

153 While many CLEF systems [38] use external general domain pre-trained
154 word-embeddings, we decided to initialize the embeddings randomly, after a
155 previous attempt initializing the embedding layer with external pre-trained
156 embeddings resulted in a performance decrease. This decay can be due to
157 the fact that external embeddings (using Word2Vec, GloVe or FastText) are
158 calculated as the result of a general learning task (the embeddings corre-
159 spond to the hidden states of the system), that might not be the kind of
160 representation needed for the task at hand. Therefore, our system obtains
161 the embeddings from scratch, recalculating them during the training process
162 of the ICD classification task itself.

163 Some of the issues we considered when we chose the method to approach
164 the ICD encoding task were: i) the large scale of the ICD classification (more
165 than 5,000 codes from around 500,000 coding items in French and Hungarian
166 and 2,500 codes from 75,000 coding items in Italian), ii) the lexical variability
167 in the texts and, iii) the need of a multilingual approach. These requirements
168 led us to approach the problem with a neural architecture, a sequence-to-
169 sequence approach which, combined with an attention mechanism, tackles
170 alignment problems, a special issue in the ICD encoding task.

171 3. Materials and methods

172 3.1. Corpora

173 We employed three datasets provided for the *Multilingual Information*
174 *Extraction* CLEF-2018 Task 1 [12]. The French dataset (CépiDc³, [40]) com-
175 prises 135,000 Death Certificates collected from 2006 to 2015. The Italian
176 dataset (ISTAT⁴) stores around 18,000 synthetic Death Certificates. They
177 were artificially built from original Death Certificates, generated in 2016,
178 with the aim of preserving confidentiality. The lines of each synthetic Death
179 Certificate were obtained from different original Certificates, always ensuring
180 topical coherence and preserving the sequence of death causes (line 1 of a
181 synthetic certificate was created using line 1 of a real certificate). The age
182 and sex of the patient were also maintained. Hence, this synthetic corpus
183 provides a realistic simulation of language and terminology found in Italian
184 Death Certificates, together with the official coding. Finally, the Hungar-
185 ian dataset consists of 100,000 Death Certificates randomly extracted from
186 a set of non-electronic Death Certificates for the year 2016 electronically
187 transcribed afterwards. This corpus was provided by the Hungarian central
188 statistical office (KSH⁵). See Table 2 for more details.

189 One of the sections in the Death Certificates [41] is a piece of text con-
190 taining one or more lines that describe the morbid diseases or conditions that
191 led to the death of an individual (an example is provided in the last column
192 of Table 3). The job of nosologists or human mortality medical coders is to

³<http://www.cepic.inserm.fr>

⁴<http://www.istat.it>

⁵<http://www.ksh.hu>

	French		Italian		Hungarian	
	Train	Test	Train	Test	Train	Test
Death Certificates	125,384	11,931	14,502	3,618	84,703	21,176
Avg. DTs per ICD10	17.59	-	7.03	-	12.74	-
#ICD10 Codes	509,103	48,948	60,955	15,789	392,020	98,264
#Uniq. ICD10s	3,723	1,806	1,443	903	3,124	2,011

Table 2: Short Description of relevant aspects of the corpus.

193 assign an ICD code to each of these conditions. In the cases where multiple
194 conditions are coded, the nosologists decide the corresponding code for each
195 DT after having seen how each condition modifies or relates to one another.
196 The dataset provided for the shared task only contains the set of lines that
197 describe those diagnoses of medical conditions. Some lines contain a single
198 diagnostic term, and its corresponding code, while others contain various
199 spontaneous DTs that correspond to a set of ICD codes.

200 All datasets were preprocessed to align the DTs at document line level as
201 expressed in the original Death Certificate, along with their corresponding
202 ICD-10 codes, as both pieces of information were contained in different files.
203 Note that a 1:1 correspondence between a DT and the ICD-10 could not be
204 assured because DTs were encoded by humans at document line level, not
205 at diagnostic term level. Therefore, only when a document line contained a
206 unique DT, then there was a 1:1 DT-ICD10 correspondence. However, when
207 a line contained several DTs, 1:1 alignment was not guaranteed (see Table 3).
208 This poses serious problems for text similarity or rule based strategies but
209 not for sequence-to-sequence systems, especially when they use attention
210 mechanisms, which allow the right alignments to be learned (see section

211 3.2.2).

ICD-10 codes	Standard Diagnostic-Term	Original Text
N179	insuffisance renale aigue	insuffisance renale aigue , masse tumorale de la tete du pancreas responsable d'une compression duodenale avec nausees / vomissements - pouvant correspondre a un 2eme primitif ou metastase
C250	tumeur maligne tete pancreas	
K566	compression duodenale	
R11	nausees vomissements	

Table 3: Example of multiple DTs within the same document line. The *Original Text* column contains several Spontaneous Diagnostic Terms. The column named *Standard Diagnostic Term* shows the standard term describing each of the multiple ICD-10 codes manually assigned to the Original Text.

212 3.2. Architecture

213 3.2.1. Baseline: Levenshtein Edit Distance

214 Edit distance is used to quantify similarity between two strings, count-
215 ing the minimum number of operations required to transform one string
216 into another. The most common metric is the Levenshtein Distance [42] in
217 where the basic edit operations are removal, insertion and substitution of
218 a single character. This metric finds the minimum distance for each spon-
219 taneous diagnostic term (SpoDT) with respect to all standard Diagnostic
220 Terms (DictDT), obtaining the best candidate match (see equation 1).

$$\text{minLev}(\text{SpoDT}, \text{DictDTs}) \quad (1)$$

221 To identify each Spontaneous DT in a text containing several DTs, we
 222 made certain assumptions. We considered colon, comma, semicolon, coordi-
 223 nation (*and or*), and certain prepositions, such as *with*, as DT separators.

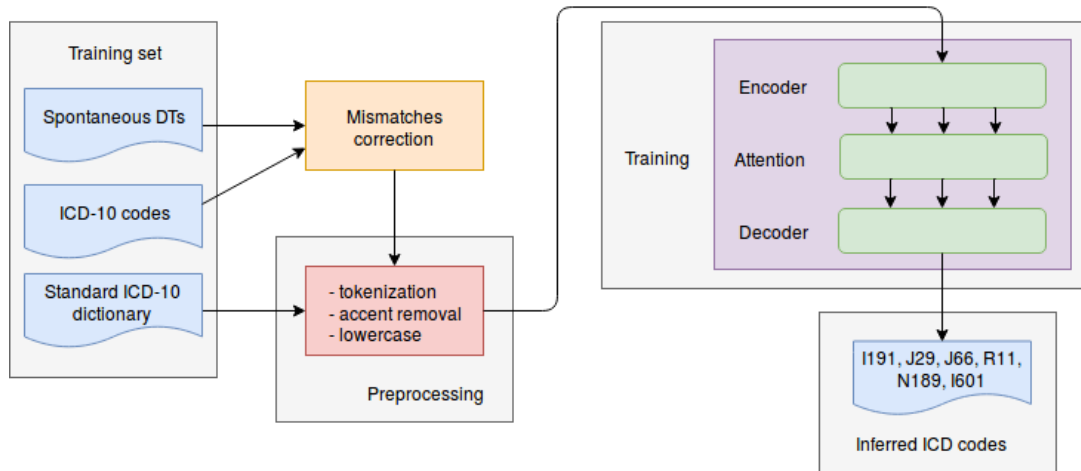


Figure 1: General architecture of the system.

224 3.2.2. ICD-10 coding as machine translation

225 In the present work we adopted a sequence-to-sequence neural machine
 226 translation (NMT) solution. Formally, having an input text $X = x_1, x_2 \dots x_n$,
 227 and an output sequence $Y = y_1 \dots y_m$, the goal is to model $P(Y|X)$. X
 228 corresponds to one document line (only those lines containing one or several
 229 DTs) and Y corresponds to one or more ICD-10 codes associated to the input
 230 DTs. Sequence-to-sequence systems calculate $P(Y|X)$ by modeling it as a
 231 step sequence where $P(Y|X) = \prod_{t=1}^m P(y_t | y_{1..t-1}, X)$. In a neural sequence-
 232 to-sequence model, the previous formula would be stated as $P(Y|X; \theta) =$
 233 $\text{softmax}(W_o s_t + b_o)$, where θ represents the neural network parameters that
 234 are automatically tuned to make the neural network minimize the error. W_o

235 corresponds to the output layer weight matrix, b_o is the bias term of this
236 layer, and s_t represents the hidden state of the neural network at step t .

237 Figure 1 presents the main components of the system. The training set
238 is composed of a set of 1:1 or N:M pairs of the form (spontaneous DT terms,
239 ICD codes). Initial preprocessing tries to minimize the mismatches in align-
240 ments of DTs and codes that appear in the relatively noisy input examples.
241 The text containing the Diagnostic Terms of the Death Certificate and the
242 corresponding codes were supplied disaggregated in different files; one file
243 stored the text containing the Diagnostic Terms as appearing in the Death
244 Certificate (Spontaneous DTs) and the location of the text occurrence (doc-
245 ument index and line index). Another file stored the corresponding ICD-10
246 codes, the standard text description of each code and the location (document
247 index and line index). In order to obtain the set of examples to train the
248 system, these two pieces of information had to be related through their lo-
249 cation in the Death Certificate. Along the process, we fixed, when possible,
250 all the location mismatches in order to improve the quality of the training
251 set. For example, there were several erroneous instances where the location
252 for an ICD code is the n^{th} line and there was no DT in the n^{th} line. In fact,
253 the spontaneous DT associated to it appeared in the $(n + 1)^{th}$ line. The
254 corrected pairs (1:1 or N:M) together with the standard ICD dictionary con-
255 taining (*standard term*, *ICD code*) pairs (1:1) were then tokenized, converted
256 to lower case and accents were removed.

257 The use of pre-trained word embeddings is generalized for many NLP
258 neural models [43, 44, 45]. Regarding NMT, when source-target data is not
259 big enough, those pre-trained embeddings have proven to be useful [46, 47].

260 The task at hand shows two differences with respect to other general NMT
261 systems. First, it is domain specific, in particular, medical domain and, sec-
262 ond, it consists in translating from words directly to codes, that is, mapping
263 between two very different spaces. We think that in this case the origin of
264 the pre-trained embeddings plays an important role on how positive their im-
265 pact in the task might be. As we mentioned in section 2, word embeddings
266 are the vectors generated in a hidden layer as a side-effect of maximizing a
267 given probability, for example the probability of predicting the center word
268 given context words (CBOW) or predicting context words given a word (Skip-
269 Gram) [48]. These tasks are general enough to make these embeddings useful
270 in many NLP tasks, but the task at hand for CLEF, namely, translate a DT
271 in an ICD, is very different in nature and too specific. The first layer of the
272 neural networks (depending on the architecture, a different neuron layer is
273 applied) is used to learn the word representations or word-embeddings. We
274 initialized it with random values, and these values were updated during the
275 training process. Therefore we can say that they were tailored for the DT-
276 ICD translation. We also experimented with pre-trained embeddings but we
277 did not obtain good results.

278 The encoder-decoder architecture using RNNs is currently the most pop-
279 ular solution for NMT, especially when combined with attention mechanisms
280 to tackle alignment problems and their limitation on long sequences by con-
281 centrating the attention on the relevant input parts. Recurrent Neural Net-
282 works or RNNs are specially useful when working with sequential data (as
283 text), thanks to its ability to maintain information about previous inputs
284 using an internal memory. RNNs carry input information across neurons by

285 means of recursive looping through an internal hidden state, thus trying to
286 maintain information about the whole sequence up to each input word. In
287 theory, RNNs can maintain information from the beginning of the sentence
288 but, in practice, this does not always happen and information from remote
289 words becomes insignificant. Long distance context is what Long Short Term
290 Memory (or LSTM) units are good at keeping inside the neural network. Us-
291 ing a LSTM unit is like adding a memory unit that can remember context
292 from the very beginning of the input.

293 Encoder-decoder architectures present two distinct blocks, the encoder
294 and the decoder. Generally speaking, an RNN encoder-decoder approach is
295 an extension of a language model. The encoder block reads the input string
296 word by word and obtains a fixed length vector representation. Then, the
297 decoder learns to sequentially predict an output code (icd_m) at position m
298 in the input string context ($DTs = w_0..w_j$) encoded by the encoder and the
299 codes predicted so far ($icd_0..icd_{m-1}$).

300 As mentioned before, and to avoid the gradient vanishing problem, the
301 RNN employed makes use of long short-term memory units (LSTM) [49].
302 When training a classic RNN using back-propagation, due to the computation
303 involved in the process, the gradients can tend to zero or infinity. With LSTM
304 units, gradients can flow unchanged. Using a gate-based system, LSTMs are
305 able to automatically regulate how much of the “previous history” context
306 should persist and how much should be renewed. Equations 2 to 6 represent
307 a basic LSTM cell as represented in Figure 2.a.

$$i_t = \sigma(W_{x_i}x_t + W_{h_i}h_{t-1} + W_{c_i}c_{t-1} + b_i) \quad (2)$$

$$\tilde{c}_t = \tanh(W_{x_c}x_t + W_{h_c}h_{t-1} + W_{c_c}c_{t-1} + b_c) \quad (3)$$

$$c_t = (1 - i_t) \odot c_{t-1} + i_t \odot \tilde{c}_t \quad (4)$$

$$o_t = \sigma(W_{x_o}x_t + W_{h_o}h_{t-1} + W_{c_o}c_t + b_o) \quad (5)$$

$$h_t = o_t \odot \tanh(c_t) \quad (6)$$

- 311 • x_i corresponds to the embedding representation of w_i .
- 312 • σ and \tanh represent the sigmoid and hyperbolic tangent, introducing
- 313 non-linearity in the network.
- 314 • t and $t - 1$ correspond to the current and previous time steps, respec-
- 315 tively.
- 316 • c_t defines the current state of the memory cell controlling how much of
- 317 the previous context $((1 - i_t) \odot c_{t-1})$ should be forgotten and how it is
- 318 updated $(i_t \odot \tilde{c}_t)$.
- 319 • i_t represents which values will get updated and \tilde{c}_t represents which new
- 320 candidates could be added.
- 321 • o_t defines, through the *sigmoid* function (σ), which part of the infor-
- 322 mation stored in the cell will go to the output.
- 323 • h_t corresponds to the hidden state. In Bi-LSTMs h_t gets calculated as
- 324 the concatenation of right to left \vec{h}_t and left to right \overleftarrow{h}_t hidden states.

325 Besides the RNN, we also experimented with other alternatives that have

326 proven to be successful for translation, as fully Convolutional Neural Net-
327 works (CNN) and self-attention (Transformer) networks.

328 A CNN network intends to identify the most relevant aspects of a Spon-
329 taneous $DT(DT = w_0..w_n)$ and represent these aspects in a fixed length
330 vector. The process consists of moving a sliding-window of k words over the
331 text obtaining several overlapping subsequences $(s_1, s_2...s_{n-k+1})$ ⁶. Then a
332 filter is applied to each subsequence with the aim of capturing some major
333 aspect of the subsequence. A filter is simply a dot-product of a given **vector**
334 representing the subsequence⁷. These weight values will be updated repeat-
335 edly over the training process using gradient descent. The result will be a
336 scalar value $(v_i = s_i.u)$. Different $u (u_1, u_2..u_j)$ filters can be applied over
337 the same subsequence s_i obtaining a vector $l_i (l_i = v_{i,1}, v_{i,2}..v_{i,j})$. Ideally
338 each $v_{i,x}$ vector will identify a different aspect of the subsequence s_i . **Next,**
339 **the “pooling” operation applies over the $l_{1:n-k+1}$ vectors to combine them**
340 **reducing the dimension to represent the initial Spontaneous DT. The reduc-**
341 **tion consists in finding the maximum value across each position in the vector**
342 **l , which indicates the most important signal in that position; this helps to**
343 **eliminate noise and also ensures that all sequences no matter how long are**
344 **represented as a fixed length vector.** This reduction is graphically represented
345 as a triangle as in Figure 2.b. The CNN decoder is similar to the encoder
346 but it has an additional attention mechanism at every layer and a fully con-
347 nected layer with a softmax to perform the actual ICD prediction. The fully

⁶Assuming that the window slides in steps of one word.

⁷Note that the representation of each subsequence is usually the concatenation of all the word vectors contained in the subsequence with a weight vector (u)

348 connected layer will use information of both the input string (one or several
 349 DTs) encoded by the encoder and the ICD codes predicted so far. For more
 350 details see [50] and [51]. Equation 7 corresponds to the generalization for
 351 a system with multiple convolutional layers (cl) of the basic convolutional
 352 filter operation.

$$h_i^{cl} = v(W^{cl}[h_{i-[k/2]}^{cl-1}; \dots; h_{i+[k/2]}^{cl-1}] + b^{cl}) + h_i^{cl-1} \quad (7)$$

353 In the Transformer architecture, a small constant step number is applied.
 354 Each step consists of a self-attention mechanism which directly models rela-
 355 tionships among all words in a source text by tuning weights. For both the
 356 encoder and the decoder the basic operation is presented in equation 8.

$$C = softmax\left(\frac{QW^Q(KW^K)^T}{\sqrt{d}}\right)VW^V \quad (8)$$

357 Q is called the Query and it is usually the last hidden state of the decoder,
 358 so it represents the target ICD codes. K is called the Key and represents the
 359 source Spontaneous DT terms (through the encoder’s last hidden state). So,
 360 the softmax, as stated, will assign bigger values to sources that are “closer”
 361 to targets, and V is referred to as the values and is equal to K , so the result of
 362 the softmax multiplied by V gives the most probable value from the source⁸.

363 In general terms, in CNN and Transformer architectures, position is not
 364 intrinsically part of the system, as opposed to RNNs⁹. Therefore, words are
 365 augmented with positional information by adding the word embedding and
 366 the positional embedding.

⁸see <http://jalamar.github.io/illustrated-transformer/> for further reading

⁹Sequentiality intrinsically encodes the position of the source items.

367 Figure 2 graphically shows the different algorithms. From the multiple
 368 toolkits available we chose Sockeye [7] because of its flexibility regarding
 369 the available architectures, including Deep Recurrent Neural Networks with
 370 Attention [52], Transformer Models with self-attention [53] and fully convo-
 371 lutional sequence-to-sequence models [50].

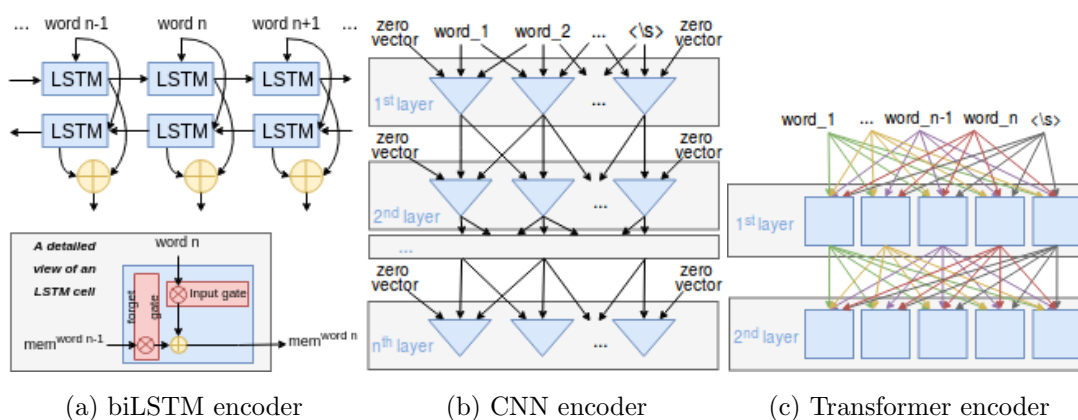


Figure 2: biLSTM, CNN and self-Transformer encoders

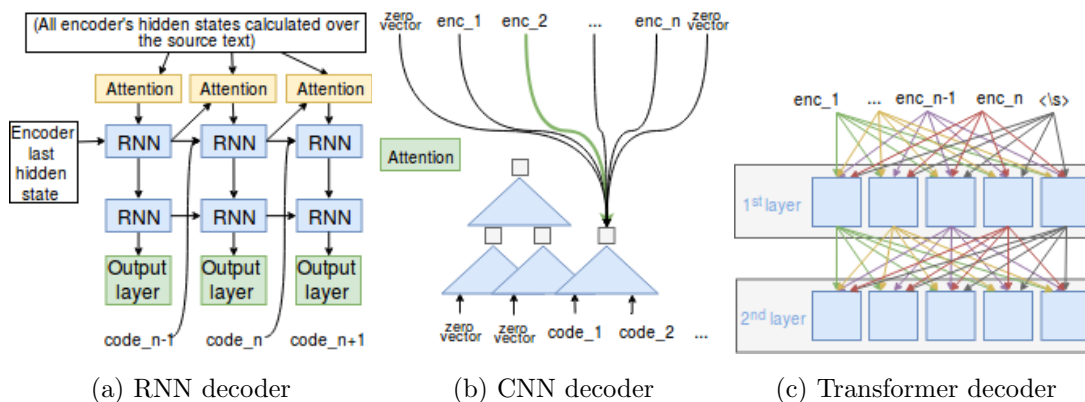


Figure 3: RNN, CNN and self-Transformer decoders

372 *3.3. Evaluation*

373 For evaluation, system performance was assessed by the usual information
374 extraction metrics: recall, precision and F-measure over the set of predicted
375 and gold standard codes (specifically, we used $\beta=1$). Matches (true positives)
376 were counted for each ICD-10 code supplied, matching the reference for the
377 associated document.

$$Precision = \frac{true\ positives}{true\ positives + false\ positives} \quad (9)$$

$$Recall = \frac{true\ positives}{true\ positives + false\ negatives} \quad (10)$$

$$F - measure = \frac{(1 + \beta^2) \times precision \times recall}{\beta \times precision \times recall} \quad (11)$$

378 **4. Results**

379 Table 4 presents the results of our system on the CLEF 2018 Shared
380 Task 1 data [12]. For the sake of comparison, we also show the results for the
381 best systems that competed in the shared task and which were described in
382 the related work section (section 2). For the French dataset, there were 18
383 official runs (or systems) from 12 teams. For the Hungarian dataset, there
384 were 9 runs from 5 teams and, for the Italian raw dataset, 12 runs from 7
385 teams.

386 In our case, we developed different models on an initial partition of the
387 data into training, validation and the initially provided test (60%-20%-20%

split). The neural network was trained on *batches* of the training set and evaluated repeatedly on the validation set until no improvement was reached. We experimented with different values for the most relevant parameters, including the embedding size for source and target tokens (256, 512 and 1024), the learning rate (from 0.0001 to 0.0005), the number of layers in encoder and decoder, the dropout rate for source and target embeddings, and the optimization algorithm. The final set of parameters uses an embedding size of 512 for source and target embeddings, 0.0003 learning rate, a single layer for encoder and decoder, no dropout and the Adam optimization algorithm. Regarding the transformer, it uses 8 heads (or parallel attention layers) in a multi-head context. We used fixed positional embeddings that are calculated deterministically using sinusoidal functions of different frequencies, related to each position [54], as they better generalize for sequences of lengths not found during training. By default, fixed positional embeddings require a size equal to that of the word embeddings. The best parameter setting was finally applied to the test set. From the built models, we selected the best two systems, which were evaluated on an unseen final test file by the shared task organizers. The final systems used the merged *train+validation* dataset in order to obtain the model parameters, leaving the initial test set for validation.

Using an Intel Xeon 3.00GHz processor with an NVIDIA TITAN X Pascal graphical processing (GPU) unit it took less than one hour for each experiment.

French			
System	Prec	Rec	F1
Baseline (Levenshtein)	0.578	0.543	0.56
Ours RNN-CNN	0.841	0.835	0.838
Ours RNN-Transformer	0.846	0.822	0.834 (-0.004)
Cossin et al. [34] dictionary-based (standard text)	0.794	0.779	0.786 (-0.052)
Cossin et al. [34] dictionary-based (ICD dict.)	0.782	0.772	0.777 (-0.061)
Gobeill and Ruch [36] instance-based learning	0.763	0.764	0.764 (-0.074)
Hungarian			
System	Prec	Rec	F1
Baseline (Levenshtein)	0.935	0.935	0.935
Ours CNN-RNN	0.970	0.955	0.963
Ours RNN-Transformer	0.968	0.954	0.961 (-0.002)
Almagro et al. [35] perceptron + One-vs-Rest	0.946	0.911	0.928 (-0.035)
Almagro et al. [35] IR	0.932	0.922	0.927 (-0.036)
Jeblee et al. [37] n-gram + word embeddings	0.922	0.897	0.910 (-0.053)
Italian			
System	Prec	Rec	F1
Baseline (Levenshtein)	0.822	0.794	0.808
Ours RNN-RNN	0.960	0.945	0.952
Ours Transformer-RNN	0.945	0.922	0.934 (-0.018)
Almagro et al. [35] perceptron + One-vs-Rest	0.917	0.875	0.895 (-0.057)
Almagro et al. [35] IR	0.931	0.861	0.895 (-0.057)
Jeblee et al. [37] n-gram + word embeddings	0.908	0.824	0.864 (-0.088)

Table 4: Performance of the best 5 systems for French, Hungarian and Italian, respectively.

411 5. Discussion

412 Table 4 shows that the sequence-to-sequence approach considerably out-
413 performs the other systems for all languages, with differences of 5.2, 3.5 and
414 5.7 absolute points in F-measure with respect to the next best system for
415 French, Hungarian and Italian, respectively.

416 Notice that the system beats the baseline as expected. Text-similarity ap-
417 proaches (Levenshtein) perform poorly on multiple DT line examples, where
418 splitting each line to find the individual DTs causes alignment problems.

419 5.1. Influence of the variability in the terms and lists of terms

420 The use of non-standard language and the variability in the DTs is a
421 source of errors (see Table 6 for the description of several error types). When
422 multiple DTs appear in the same line, besides language variability, alignment
423 issues arise, because identifying individual DTs becomes tricky. For exam-
424 ple, two DTs separated by a comma, as in “*evolution terminale , insuffisance*
425 *cardiaque*” (ICD-codes R999 I509), might not obtain the same code as when
426 they appear with a preposition, “*evolution terminale d’une insuffisance car-*
427 *diaque*” (ICD-code I509). However, having multiple terms can be helpful
428 sometimes. Terms might be interrelated, so a term might thus help to find
429 the right code for another one (see Figure 6).

430 5.2. Influence of neural architecture

431 Although we experimented with other combinations of encoder-decoder
432 pairs, Table 4 presents the results of our best two systems for each language.
433 The table shows important differences depending on the neural network ar-
434 chitecture employed. There is no unique encoder-decoder combination that

435 performs best for all languages, the best systems being RNN-CNN, CNN-
436 RNN and RNN-RNN for French, Hungarian and Italian, respectively. Re-
437 garding the decoders, RNNs seem to work better for Italian, as the best
438 combinations use an RNN decoder. In Hungarian the second best combi-
439 nation uses an RNN as decoder and, finally, none of the best combinations
440 for French employs an RNN. The most plausible explanation for this fact
441 originates from the data itself. RNNs are intrinsically sequential and conse-
442 quently the order of the ICDs might be an issue in lines with multiple ICDs.
443 It turns out that in Italian only 2% of the data corresponds to multiple ICD
444 lines with an average length of 2.4 ICDs, while in Hungarian it is a 3% and
445 in French it is around a 5% with an average length of 2.9 both. From those
446 multiple ICD lines, Italian is the language showing the lowest order variabil-
447 ity, then Hungarian and, finally, French with the highest order variability.
448 As for the encoder, the opposite applies, because the language showing the
449 highest word order freedom is Hungarian, then Italian and finally French.
450 This suggests that the optimal algorithms should be selected after a careful
451 experimentation.

452 Besides the general architecture, we also tested variations of the different
453 hyperparameters. Our main experience is that the difference comes mostly
454 from varying the architectures (CNN, RNN and transformer), rather than
455 from adjusting the hyperparameters, with minor variations from the standard
456 default values.

457 We made preliminary experiments using subword units known as byte-pair
458 encodings [6] that try to define smaller segments than the word itself. This
459 is based on the intuition that various word classes are translatable via units

460 smaller than words, for instance names (via character copying or transliter-
461 ation), compounds (via compositional translation), and cognates and loan-
462 words (via phonological and morphological transformations). Although this
463 can in principle be useful for dealing with unknown words, we did not find
464 any significant improvement.

465 5.3. Interpretability

466 The examples in Table 3 show that the connection between the codes and
467 the parts of the DTs is not annotated, obscuring the reasons behind specific
468 coding decisions. As Li et al. pointed out [55]: “unlike traditional feature-
469 based classifiers that assign and optimize weights to varieties of human in-
470 terpretable features (parts-of-speech, named entities, word shapes, syntactic
471 parse features etc) the behavior of deep learning models is much less eas-
472 ily interpreted”. Recent works have made an effort towards this direction
473 [56, 57].

474 Figure 4 shows how the network processes the DT input, using a vi-
475 sualization tool for neural sequence-to-sequence models [58]. The different
476 values calculated in the process get represented graphically allowing [users](#) to
477 understand the decisions at each step. The figure presents the three main
478 components of the system: the encoder (in blue), the attention (linking en-
479 coder and decoder) and the decoder (in yellow). For instance, in Figure 4.a,
480 for the input *hepatite c cardiopathie rythmique*, the right codes are correctly
481 predicted as *B182* (hepatite C) and *I499* (Cardiopathie rythmique). The
482 width of the attention lines represents the weight of each encoder state (the
483 input) over the decoder predictions (the output). For example, the word “c”
484 plays a relevant role to decide the right code for “hepatite c” (Figure 4.a),

485 as all the codes at the top k candidates refer to hepatic diseases. The can-
486 didates, ranked by their probability, are B182 (*Hépatite virale chronique C*),
487 K759 (*Maladie inflammatoire du foie, sans précision*), B189 (*Hépatite virale*
488 *chronique, sans précision*) and K746 (*Cirrhoses, autres et sans précision*).
489 The attention mechanism is one of the keys that helped to obtain accurate
490 predictions, ~~comparing~~ compared to other participant groups that did not
491 make use of attention. The lower part of Figure 4 shows the alternative
492 paths that are evaluated when searching for the best code assignment, rep-
493 resenting the path likelihood by the width of the connecting lines.

494 Figure 5 shows the encoder contexts analogous to the word “aomi” (Fig-
495 ure 4.b), presenting DTs related to the input. This way, the human coder can
496 inspect similar encoding states which in this case allow to ascertain the link-
497 ing between “aomi” and “arteriopathie oblitérante des membres inférieurs”
498 (as an acronym).

499 Figure 6 presents the alignment examples of two DTs and ICD-10 codes,
500 showing the result of the sequence-to-sequence alignment. In addition to
501 assigning codes with good accuracy, our system provides an interpretable
502 result, aligning each code with its corresponding piece of text as in [59].
503 This information is not rendered by most ML approaches which, although
504 accurate, do not provide any helpful information besides the result itself.
505 This is relevant in medical environments where hospitals/clinicians **require**
506 an understandable way to find the most informative evidence. We provided
507 a physician with a sample of 78 ICDs appearing in 15 lines, with and with-
508 out the alignment matrices. Overall, the physician found the system very
509 helpful. Finding the right ICD for a given DT from scratch, that is, with-

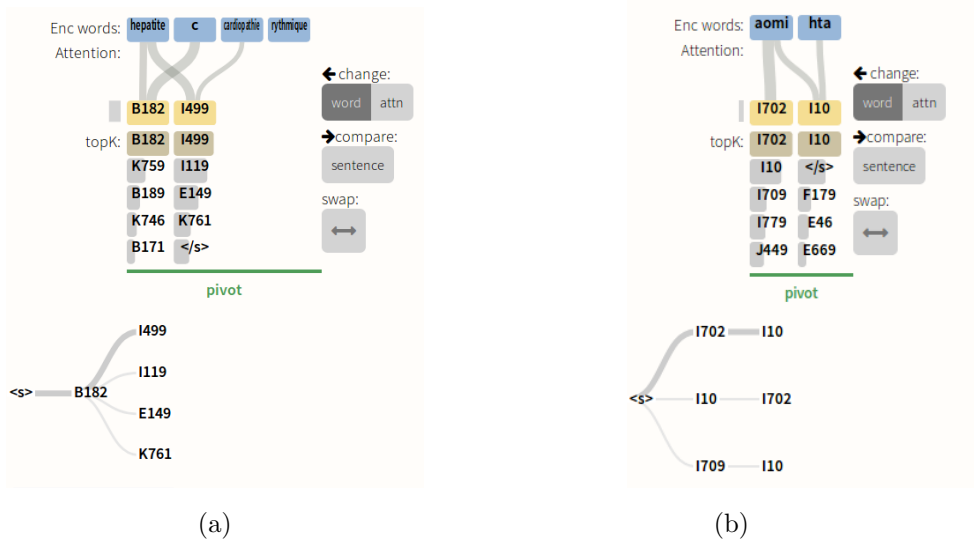


Figure 4: Examples of the ICD code generation process.

510 out having any system prediction at all, ~~implied to start looking~~ required
 511 the human annotator to look up the DT in the ICD coding classification as
 512 provided by the WHO (World Health Organization). This was hard since
 513 DTs do not usually match the standard terms, so when there was no per-
 514 fect match the expert had to locate in the hierarchy the right chapter (DT
 515 generalization, e.g. *insuffisance renale terminale* corresponds to chapter XIV
 516 codes from N00N99, for Maladies de l'appareil génito-urinaire) and then dig
 517 ~~up~~ for the right code. Having a prediction flips the process: the physician
 518 started by looking up ~~for~~ the predicted ICD code to obtain the standard
 519 term associated to it, and if the standard term was a synonym of the DT to
 520 be classified, the process finished successfully. But even when the prediction
 521 was not correct, as the ICDs follow a hierarchical structure, a partial predic-
 522 tion narrowed the search space reducing the coding time (e.g., *insuffisance*

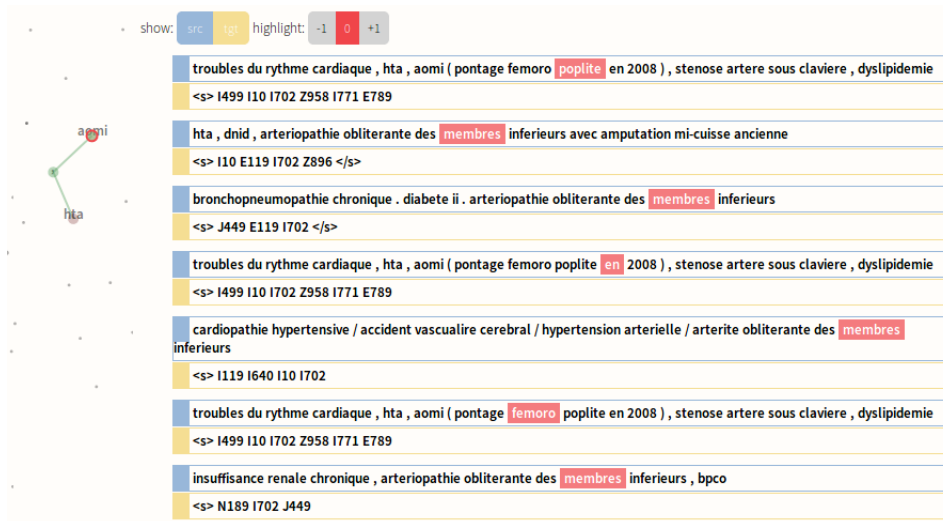


Figure 5: Related encoder contexts of the word “aomi” (Figure 4.b), showing DTs related to the input (e.g. “arteriopathie oblitérante des membres inférieurs”).

523 *renale terminale* was assigned the ICD code N180 which, although incorrect,
 524 is close to the correct code, N185). The alignment matrices, in particular,
 525 were useful when errors occur, ~~specially~~ especially in long lines that contain
 526 several DTs, as the alignments showed a one-to-one correspondence between
 527 a term and the predicted ICD, allowing the physician to focus on that exact
 528 pair and, if the prediction was wrong, the physician did not need to look for
 529 any other DT that could match the ICD.

530 In figure 6.a, we see how the DT *acido-cetose diabetique*, containing a
 531 spelling error, is paired with the code E141 that corresponds to the standard
 532 term *acidocetose diabetique*. Figure 6.b depicts how the system finds different
 533 types of evidence assigned to each DT word with respect to an ICD-10 code,
 534 illustrated by the DT *insuffisance respiratoire restrictive*, where each word
 535 adds value to the alignment with the term E274, and how some words are

536 more relevant than others.

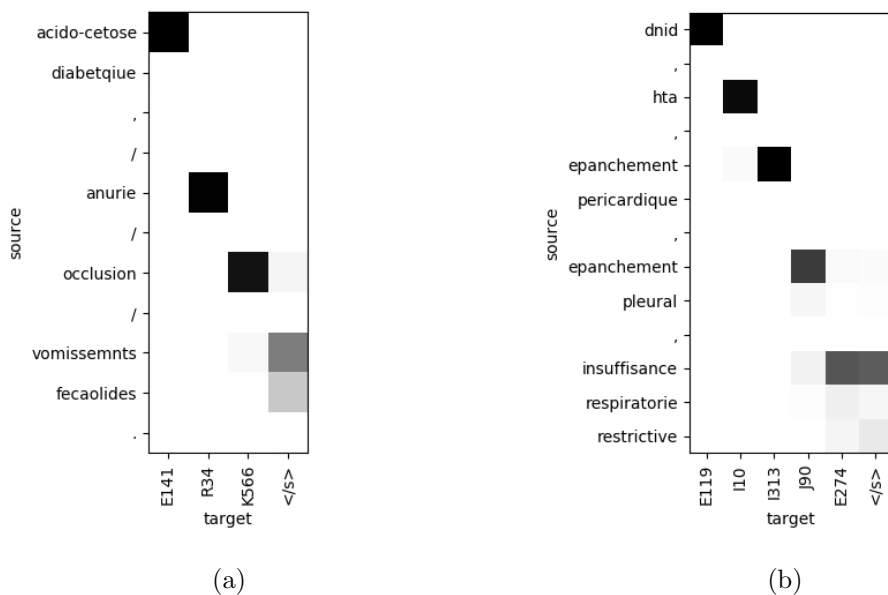


Figure 6: Alignment examples.

537 Additionally, word embeddings calculated in the translation process cap-
 538 ture similarities between elements, such as semantically similar equivalents,
 539 variants or spelling errors. Table 5 presents the closest words in the embed-
 540 ding space to a set of selected words. For example, given the French word
 541 *bronchique*, the set of close word embeddings includes variants (*bronchiques*),
 542 spelling errors (*brochique*, *bronchique*), or semantically related terms (*tra-*
 543 *cheobronchique*).

544 5.4. Effect of algorithms and training corpus size

545 We found it interesting to explore the relation between the types of algo-
 546 rithms and certain characteristics of the corpora, such as size and number of

Term (English)	French	Italian
bronchial	bronchique	bronchit
	brochique	bronchite
	bronchique	broncopat
	bronchiques	broncopneumopat
	bronchioque	tracheobronchite
	tracheobronchique	broncopneumopatia
dehydration	deshydratation	disidratazione
	hypovolemie	ipovolemia
	deshydrataion	idratazione
	deshydratee	deidratato
	dehydration	disidratativa

Table 5: Example of close embeddings for the terms “bronchial” and “dehydration” in French and Italian.

547 codes. To do so, we artificially built several ~~size~~ differently-sized training cor-
548 pora. Figure 7 shows how Transformer for encoding and CNN for decoding
549 present the best results for small datasets, with an F-score difference of ~~10~~
550 ~~0.10 absolute points with~~ from other approaches. However, it also shows the
551 worst performance on the full dataset. Conversely, architectures based on
552 CNN and RNN do not obtain a good accuracy with small training sizes, but
553 achieve the best results with bigger datasets. This could be due to the fact
554 that Transformer is less constrained to positional information, generalizes
555 better, and suffers less from sparsity. The attentional Transformer-RNN pair
556 is the most regular architecture for both small and big datasets. These re-
557 sults present interesting ideas for the implementation of neural architectures,
558 depending on the language or the amount of data.

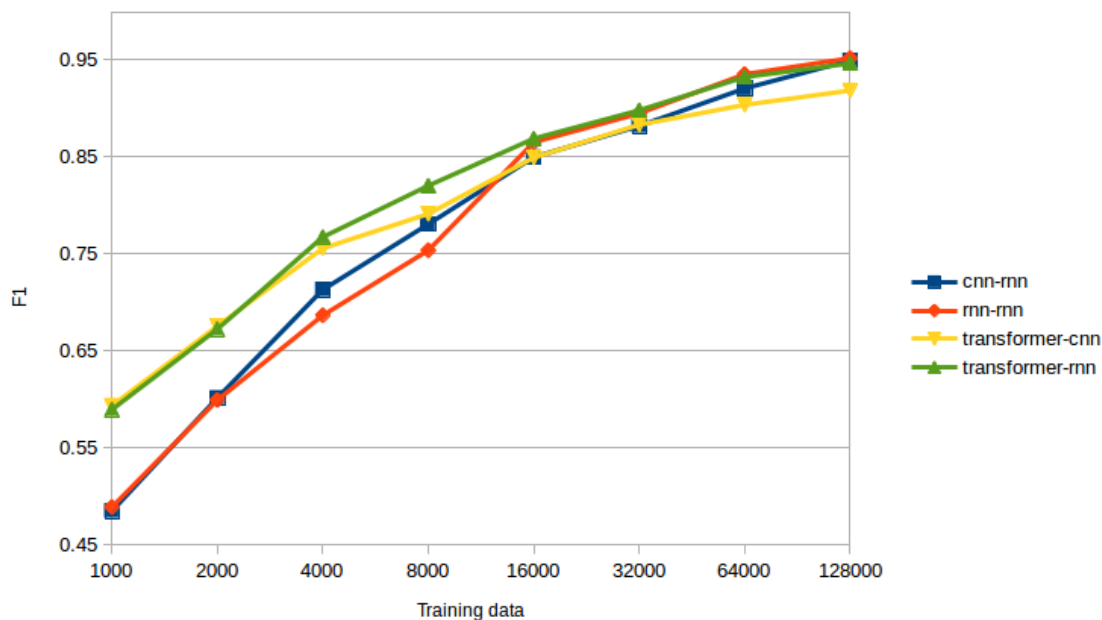


Figure 7: Effect of increasing the training size (Italian) for different encoder and decoders.

559 *5.5. Error analysis*

560 In order to understand the results and with the aim of improving them in
 561 future developments, we analyzed the errors manually. Table 6 exemplifies
 562 several error sources. This analysis was divided into two groups: i) document
 563 lines with a single DT and ii) lines with multiple DTs. These are a general
 564 source of error in both cases:

- 565 • Abbreviation: the standard form appears abbreviated and so the in-
 566 formation is misunderstood. In some cases it carries an incorrect code
 567 and in others the addition of a new code.
- 568 • Spelling error: the spontaneous DT is misspelled and so the system
 569 obtains an incorrect code.

- 570 • Superclass: the system gives the spontaneous DT the code correspond-
571 ing to its direct parent in the ICD-10-CM hierarchy (that means a more
572 general code).
- 573 • Information inclusion/exclusion: the system does not identify part of
574 the DT, so it generates a code mismatch. The same thing happens
575 whenever the system considers more tokens for the code assignment.

576 However, other types of problems, like the omission of commas, or the
577 variations in the use of coordination or certain prepositions only appear when
578 DT lists are used, as this may cause alignment problems. In fact, the system
579 performs better with the one-to-one (1:1) cases which are not prompted to
580 alignment issues. This difference, as expected, was even bigger when using
581 Levenshtein.

582 **6. Conclusion**

583 This work tackles medical record classification following the ICD-10 stan-
584 dard in a multilingual setting. The classification problem is hard for several
585 reasons: 1) the gap between spontaneous and standard language; 2) a large-
586 scale classification task, with thousands of possible classes and, 3) in real
587 data, in most cases, there is no 1:1 alignment between DTs and ICD-10
588 codes.

589 We present and evaluate different neural network architectures for multi-
590 class document classification as a sequence-to-sequence problem. The system
591 is also able to learn to identify highly-predictive locations for each label,
592 providing satisfactory explanations for its predictions. The system is also able

Error analysis with examples	
1:1	
Abbreviation	G: oap asphyxique I501 S: oedeme aigu pulmonaire asphyxique I501 R090
Multiword Spelling error	G: arret cardio-circulatoire I469 G: detresse respiratoire J960 S: detresse respiratoiore R092
Superclass	G: acfa I489 S: acfa I48
Info. missing	G: plaie cardiaque operatoire Y600 S: plaie cardiaque S269
Info. inserted	G: opereee Z924 S: neoplasie analse opereee C169 Z924
Prep. inserted	G: syndrome glissement R453 S: syndrome de glissement R54
...	
1:N	
All the cases in 1:1 and in addition the following:	
Term unification by comma omission	G: deshydratation , demence type alzheimer , tumeur sein E86 G309 D486 S: deshydratation demence type alzheimer tumeur du sein E86 G309 C509
Term unification by comma substitution	G: hemorragie encephalique , hemorragie ventriculaire I619 I615 S: hemorragie encephalique et ventriculaire S062
...	

Table 6: Error cases detected (French). G refers to the standard form and S to the spontaneous form. The sequence of one or more DTs is followed by the corresponding ICD codes, either manual or automatic.

593 to deal with real non-aligned data which is difficult for some other approaches,
594 such as text similarity based models.

595 Our best model showed high-quality results, establishing a new state-of-
596 the-art, and this fact opens a promising avenue for the task of automatically
597 assigning ICD-10 codes to medical documents. The method is language inde-
598 pendent, allowing efficient training, given only a set of annotated documents,
599 and does not require complex feature engineering.

600 **Acknowledgements**

601 This work was supported by the Spanish Ministry of Science and In-
602 novation (TUNER: TIN2015-65308-C5-1-R, PROSA-MED: TIN2016-77820-
603 C3-1-R, DOMINO: PGC2018-102041-B-I00), and the Basque Government
604 (MODENA: KK-2018/00087).

605 We gratefully acknowledge the support of NVIDIA Corporation with the
606 donation of the Titan X Pascal GPU used for this research.

607 **References**

- 608 [1] D. Lang, Consultant Report-Natural Language Processing in the Health
609 Care Industry, Cincinnati Children’s Hospital Medical Center, Winter
610 (2007).
- 611 [2] W. H. Organization, et al., International Classification of Diseases ICD
612 10 Classification of Mental and Behavioural Disorders: Clinical Descrip-
613 tions and Diagnostic Guidelines, World Health Organisation, 1993.

- 614 [3] K. Cho, B. van Merriënboer, G. Gulcehre, F. Bougares, H. Schwenk,
615 Y. Bengio, Learning Phrase Representations using RNN Encoder-
616 Decoder for Statistical Machine Translation, in: Conference on Em-
617 pirical Methods in Natural Language Processing (EMNLP) (2014).
- 618 [4] I. Sutskever, O. Vinyals, Q. Le, Sequence to Sequence Learning with
619 Neural Networks, in: Advances in Neural Information Processing Sys-
620 tems 27: Annual Conference on Neural Information Processing Systems
621 (NIPS), pp. 3104–3112 (2014).
- 622 [5] G. Klein, Y. Kim, Y. Deng, J. Senellart, A. Rush, OpenNMT: Open-
623 Source Toolkit for Neural Machine Translation, in: Proceedings of ACL
624 2017, System Demonstrations, Association for Computational Linguis-
625 tics, Vancouver, Canada, 2017, pp. 67–72.
- 626 [6] R. Sennrich, O. Firat, K. Cho, A. Birch, B. Haddow, J. Hitschler,
627 M. Junczys-Dowmunt, S. Läubli, A. V. Miceli Barone, J. Mokry,
628 M. Nadejde, Nematus: a Toolkit for Neural Machine Translation, in:
629 Proceedings of the Software Demonstrations of the 15th Conference of
630 the European Chapter of the Association for Computational Linguis-
631 tics, Association for Computational Linguistics, Valencia, Spain, 2017,
632 pp. 65–68.
- 633 [7] F. Hieber, T. Domhan, M. Denkowski, D. Vilar, A. Sokolov, A. Clifton,
634 M. Post, The SOCKEYE Neural Machine Translation Toolkit at AMTA
635 2018, in: Proceedings of the 13th Conference of the Association for
636 Machine Translation in the Americas (Volume 1: Research Papers),

- 637 Association for Machine Translation in the Americas, 2018, pp. 200–
638 207.
- 639 [8] J. P. Pestian, C. Brew, P. Matykiewicz, D. J. Hovermale, N. Johnson,
640 K. B. Cohen, W. Duch, A Shared Task Involving Multi-label Classifi-
641 cation of Clinical Free Text, in: Workshop on BioNLP’07: Biological,
642 Translational, and Clinical Language Processing, 2007.
- 643 [9] Ö. Uzuner, B. R. South, S. Shen, S. L. DuVall, 2010 i2b2/VA challenge
644 on concepts, assertions, and relations in clinical text, JAMIA 18 (2011)
645 552–556.
- 646 [10] R. Bossy, W. Golik, Z. Ratkovic, D. Valsamou, P. Bessières, C. Nédellec,
647 Overview of the gene regulation network and the bacteria biotope tasks
648 in BioNLP’13 shared task, BMC Bioinformatics 16 (2015) S1.
- 649 [11] A. Névéol, R. N. Anderson, K. B. Cohen, C. Grouin, T. Lavergne,
650 G. Rey, A. Robert, C. Rondet, P. Zweigenbaum, CLEF eHealth 2017
651 Multilingual Information Extraction task overview: ICD10 coding of
652 death certificates in English and French, in: CLEF 2017 Evaluation
653 Labs and Workshop: Online Working Notes, CEUR-WS.
- 654 [12] A. Névéol, A. Robert, F. Grippio, C. Morgand, C. Orsi, L. Pelikán,
655 L. Ramadier, G. Rey, P. Zweigenbaum, CLEF eHealth 2018 Multilin-
656 gual Information Extraction task Overview: ICD10 Coding of Death
657 Certificates in French, Hungarian and Italian, in: CLEF 2018 Evalua-
658 tion Labs and Workshop: Online Working Notes, CEUR-WS, September
659 (2018).

- 660 [13] G. Schmidt, ICD-10 coding using machine learning, <http://www.gregoryschmidt.ca/writing/icd-coding>, 2016.
- 661
- 662 [14] S. V. Pakhomov, J. D. Buntrock, C. G. Chute, Automating the Assign-
663 ment of Diagnosis Codes to Patient Encounters Using Example-based
664 and Machine Learning Techniques, *Journal of the American Medical
665 Informatics Association* 13 (2006) 516–525.
- 666 [15] S. Boytcheva, I. Nikolova, E. Paskaleva, G. Angelova, D. Tcharaktchiev,
667 N. Dimitrova, Structuring of status descriptions in hospital patient
668 records, in: the Proceedings 2nd International Workshop on Building
669 and Evaluating Resources for BioMedical Text Mining, associated to
670 the 7th Int. Conf. on Language Resources and Evaluation (LREC-2010),
671 ELRA, pp. 31–36.
- 672 [16] G. Karystianis, A. J. Nevado, C.-H. Kim, A. Dehghan, J. A. Keane,
673 G. Nenadic, Automatic mining of symptom severity from psychiatric
674 evaluation notes, in: *International journal of methods in psychiatric
675 research*.
- 676 [17] R. Farkas, G. Szarvas, Automatic construction of rule-based ICD-9-CM
677 coding systems, *BMC bioinformatics* 9 (2008) 1–9.
- 678 [18] A. Pérez, A. Atutxa, A. Casillas, S. A. Gojenola, K., Inferred joint
679 multigram models for medical term normalization according to ICD,
680 *International Journal of Medical Informatics*, **110**, pp. 111–117 (2018).
- 681 [19] R. A. Kavuluru A, L. Y., An empirical evaluation of supervised learning

- 682 approaches in assigning diagnosis codes to electronic medical records,
683 *Artif Intell Med.* 65 (2015) 155–166.
- 684 [20] M. Almagro, R. Martínez, V. Fresno, S. Montalvo, Estudio prelimi-
685 nar de la anotación automática de códigos CIE-10 en informes de alta
686 hospitalarios, *Procesamiento del Lenguaje Natural* 60 (2018) 45–52.
- 687 [21] B. Koopman, G. Zuccon, A. Nguyen, A. Bergheim, N. Grayson, Auto-
688 matic ICD-10 classification of cancers from free-text death certificates,
689 *International Journal of Medical Informatics* 84 (2015) 956 – 965.
- 690 [22] B. Koopman, S. Karimi, A. Nguyen, R. McGuire, D. Muscatello,
691 M. Kemp, D. Truran, M. Zhang, S. Thackway, Automatic classification
692 of diseases from free-text death certificates for real-time surveillance,
693 *BMC Medical Informatics and Decision Making* 15 (2015) 53.
- 694 [23] K. Xu, M. Lam, J. Pang, X. Gao, C. Band, P. Mathur, F. Papay, A. K.
695 Khanna, J. B. Cywinski, K. Maheshwari, P. Xie, E. Xing, Multimodal
696 Machine Learning for Automated ICD Coding, *CoRR abs/1810.13348*
697 (2018).
- 698 [24] T. Young, D. Hazarika, S. Poria, E. Cambria, Recent Trends in Deep
699 Learning Based Natural Language Processing, *CoRR abs/1708.02709*
700 (2017).
- 701 [25] R. Miotto, F. Wang, S. Wang, X. Jiang, J. T. Dudley, Deep learn-
702 ing for healthcare: review, opportunities and challenges, *Briefings in*
703 *bioinformatics* 19 (2017) 1236–1246.

- 704 [26] J. Mullenbach, S. Wiegreffe, J. Duke, J. Sun, J. Eisenstein, Explainable
705 Prediction of Medical Codes from Clinical Text, in: Proceedings of
706 the 2018 Conference of the North American Chapter of the Association
707 for Computational Linguistics: Human Language Technologies, Volume
708 1 (Long Papers), Association for Computational Linguistics, 2018, pp.
709 1101–1111.
- 710 [27] H. Shi, P. Xie, Z. Hu, M. Zhang, E. P. Xing, Towards automated ICD
711 coding using deep learning, arXiv preprint arXiv:1711.04075 (2017).
- 712 [28] S. J. Pan, Q. Yang, A survey on transfer learning, IEEE Transactions
713 on knowledge and data engineering 22 (2010) 1345–1359.
- 714 [29] M. Zeng, M. Li, Z. Fei, Y. Yu, Y. Pan, J. Wang, Automatic ICD-9
715 coding via deep transfer learning, Neurocomputing 324 (2019) 43–50.
- 716 [30] K. Liu, S. Peng, J. Wu, C. Zhai, H. Mamitsuka, S. Zhu, MeSHLabeler:
717 improving the accuracy of large-scale MeSH indexing by integrating di-
718 verse evidence, Bioinformatics 31 (2015) i339–i347.
- 719 [31] M. Ebersbach, R. Herms, M. Eibl, Fusion Methods for ICD10 Code
720 Classification of Death Certificates in Multilingual Corpora, in: CLEF
721 2017 Conference and Labs of the Evaluation Forum, Online Working
722 Notes, CEUR-WS, September (2017).
- 723 [32] P. Zweigenbaum, T. Lavergne, Multiple Methods for Multi-class, Multi-
724 label ICD-10 Coding of Multi-granularity, Multilingual Death Certifi-
725 cates, in: CLEF 2017 Conference and Labs of the Evaluation Forum,
726 Online Working Notes, CEUR-WS, September (2017).

- 727 [33] Z. Miftahutdinov, E. Tutubalina, KFU at CLEF eHealth 2017 Task
728 1: ICD-10 Coding of English Death Certificates with Recurrent Neural
729 Networks, in: CLEF 2017 Conference and Labs of the Evaluation Forum,
730 Online Working Notes, CEUR-WS, September (2017).
- 731 [34] S. Cossin, V. Jouhet, F. Mouglin, G. Diallo, F. Thiessard, IAM at
732 CLEF eHealth 2018 : Concept Annotation and Coding in French Death
733 Certificates, in: CLEF 2018 Evaluation Labs and Workshop: Online
734 Working Notes, CEUR-WS, September (2018).
- 735 [35] M. Almagro, S. Montalvo, A. Diaz de Ilarraza, A. Pérez, MAMTRA-
736 MED at CLEF eHealth 2018: A Combination of Information Retrieval
737 Techniques and Neural Networks for ICD-10 Coding of Death Certifi-
738 cates, in: CLEF 2018 Evaluation Labs and Workshop: Online Working
739 Notes, CEUR-WS, September (2018).
- 740 [36] J. Gobeill, P. Ruch, Instance-based Learning for ICD10 Categorization,
741 in: CLEF 2018 Evaluation Labs and Workshop: Online Working Notes,
742 CEUR-WS, September (2018).
- 743 [37] S. Jeblee, A. Budhkar, S. Milić, J. Pinto, C. Pou-Prom, K. Vishnubhotla,
744 G. Hirst, F. Rudzicz, Toronto CL at the CLEF 2018 eHealth Challenge
745 Task 1, in: CLEF 2018 Evaluation Labs and Workshop: Online Working
746 Notes, CEUR-WS, September (2018).
- 747 [38] J. Ševa, M. Sängler, U. Leser, WBI at CLEF eHealth 2018 Task 1:
748 Language-independent ICD-10 Coding using Multi-lingual Embeddings

- 749 and Recurrent Neural Networks, in: CLEF 2018 Evaluation Labs and
750 Workshop: Online Working Notes, CEUR-WS, September (2018).
- 751 [39] F. Duarte, B. Martins, C. S. Pinto, M. J. Silva, Deep neural models
752 for ICD-10 coding of death certificates and autopsy reports in free-text,
753 Journal of biomedical informatics 80 (2018) 64–77.
- 754 [40] T. Lavergne, A. Névéol, A. Robert, C. Grouin, G. Rey, P. Zweigenbaum,
755 Dataset for ICD-10 Coding of Death Certificates: Creation and Usage,
756 in: Proceedings of the Fifth Workshop on Building and Evaluating Re-
757 sources for Biomedical Text Mining (BioTxtM 2016).
- 758 [41] G. Pavillon, F. Laurent, Certification et codification des causes
759 médicales de décès, Bull Epidemiol Hebd 30 (2003) 134–138.
- 760 [42] V. I. Levenshtein, Binary codes capable of correcting deletions, inser-
761 tions, and reversals, in: Soviet physics doklady, volume 10, pp. 707–710.
- 762 [43] G. Lample, M. Ballesteros, S. Subramanian, K. Kawakami, C. Dyer,
763 Neural architectures for named entity recognition, arXiv preprint
764 arXiv:1603.01360 (2016).
- 765 [44] X. Ma, E. Hovy, End-to-end sequence labeling via bi-directional lstm-
766 cnn-crf, arXiv preprint arXiv:1603.01354 (2016).
- 767 [45] Y. Kim, Convolutional neural networks for sentence classification, arXiv
768 preprint arXiv:1408.5882 (2014).
- 769 [46] M. Artetxe, G. Labaka, E. Agirre, K. Cho, Unsupervised neural machine
770 translation, arXiv preprint arXiv:1710.11041 (2017).

- 771 [47] M. Neishi, J. Sakuma, S. Tohda, S. Ishiwatari, N. Yoshinaga, M. Toy-
772 oda, A bag of useful tricks for practical neural machine translation:
773 Embedding layer initialization and large batch size, in: Proceedings of
774 the 4th Workshop on Asian Translation (WAT2017), pp. 99–109.
- 775 [48] T. Mikolov, K. Chen, G. S. Corrado, J. Dean, Efficient Estimation of
776 Word Representations in Vector Space, CoRR abs/1301.3781 (2013).
- 777 [49] S. Hochreiter, J. Schmidhuber, Long Short-Term Memory, Neural com-
778 putation 9 (1997) 1735–1780.
- 779 [50] J. Gehring, M. Auli, D. Grangier, D. Yarats, Y. Dauphin, Convolutional
780 Sequence to Sequence Learning, in: Proceedings of the 34th Interna-
781 tional Conference on Machine Learning.
- 782 [51] J. Gehring, M. Auli, A novel approach to neural ma-
783 chine translation, [https://code.fb.com/ml-applications/
784 a-novel-approach-to-neural-machine-translation/](https://code.fb.com/ml-applications/a-novel-approach-to-neural-machine-translation/), 2017.
- 785 [52] D. Bahdanau, K. Cho, Y. Bengio, Neural Machine Translation by Jointly
786 Learning to Align and Translate, in: arXiv, arXiv:1409.0473.
- 787 [53] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez,
788 L. Kaiser, I. Polosukhin, Attention Is All You Need, in: Advances in
789 Neural Information Processing Systems, pp. 6000–6010.
- 790 [54] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez,
791 L. u. Kaiser, I. Polosukhin, Attention is all you need, in: I. Guyon, U. V.
792 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett

- 793 (Eds.), *Advances in Neural Information Processing Systems 30*, Curran
794 Associates, Inc., 2017, pp. 5998–6008.
- 795 [55] J. Li, X. Chen, E. Hovy, D. Jurafsky, Visualizing and Understand-
796 ing Neural Models in NLP, in: *Proceedings of the 2016 Conference of*
797 *the North American Chapter of the Association for Computational Lin-*
798 *guistics: Human Language Technologies*, Association for Computational
799 Linguistics, 2016, pp. 681–691.
- 800 [56] Y. Ding, Y. Liu, H. Luan, M. Sun, Visualizing and Understanding Neu-
801 ral Machine Translation, in: *Proceedings of the 55th Annual Meeting of*
802 *the Association for Computational Linguistics (Volume 1: Long Papers)*,
803 Association for Computational Linguistics, 2017, pp. 1150–1159.
- 804 [57] B. C. Kwon, M. Choi, J. T. Kim, E. Choi, Y. B. Kim, S. Kwon, J. Sun,
805 J. Choo, RetainVis: Visual Analytics with Interpretable and Interac-
806 tive Recurrent Neural Networks on Electronic Medical Records, *CoRR*
807 [abs/1805.10724](https://arxiv.org/abs/1805.10724) (2018).
- 808 [58] H. Strobelt, S. Gehrmann, M. Behrisch, A. Perer, H. Pfister, A. Rush,
809 Seq2Seq-Vis: A Visual Debugging Tool for Sequence-to-Sequence Mod-
810 els, *ArXiv e-prints* (2018).
- 811 [59] I. Lopez-Gazpio, M. Maritxalar, A. Gonzalez-Agirre, G. Rigau, L. Uria,
812 E. Agirre, Interpretable semantic textual similarity: Finding and ex-
813 plaining differences between sentences, *Knowledge-Based Systems* 119
814 (2017) 186 – 199.

815 **7. Summary Points**

816 What was already known on the topic?

- 817 • The problem presents a complex characterization due to non-standard
818 language variation, spontaneous writing, large-scale multi-class classi-
819 fication or DT-ICD10 alignment issues.
- 820 • There are varied approaches, ranging from knowledge-based solutions
821 to statistical and deep learning ones.
- 822 • Most Machine Learning approaches, although accurate, do not offer
823 any helpful clue about the encoding decision besides the result itself.

824 What does this work add?

- 825 • Sequence-to-sequence deep learning approaches outperform other sys-
826 tems by a considerable margin for all languages.
- 827 • We have performed an exhaustive study of different sequence-to-sequence
828 architectures, showing that there is no unique encoder-decoder com-
829 bination that performs best for all languages, as we show important
830 differences with respect to the neural network architecture employed.
- 831 • Apart from assigning the codes with good accuracy, the system also pro-
832 vides an interpretable result, aligning each code with its corresponding
833 piece of text.