

# A Section Identification Tool: towards HL7 CDA/CCR Standardization in Spanish Discharge Summaries

Iakes Goenaga<sup>a</sup>, Xabier Lahuerta<sup>a</sup>, Aitziber Atutxa<sup>b</sup>, Koldo Gojenola<sup>b</sup>

*HiTZ Basque Center for Language Technology*

*<http://www.hitz.eus>*

*University of the Basque Country (UPV/EHU), Spain*

<sup>a</sup>*Faculty of Computer Science, P<sup>o</sup> Manuel Lardizabal, 1 — 20018 Donostia-San Sebastián*

<sup>b</sup>*School of Engineering, Paseo Rafael Moreno Pitxitxi, 3 — 48013 Bilbao*

---

## Abstract

**Background.** Nowadays, with the digitalization of healthcare systems, huge amounts of clinical narratives are available. However, despite the wealth of information contained in them, interoperability and extraction of relevant information from documents remains a challenge.

**Objective.** This work presents an approach towards automatically standardizing Spanish Electronic Discharge Summaries (EDS) following the HL7 Clinical Document Architecture. We address the task of section annotation in EDSs written in Spanish, experimenting with three different approaches, with the aim of boosting interoperability across healthcare systems and hospitals.

**Methods.** The paper presents three different methods, ranging from a knowledge-based solution by means of manually constructed rules to super-

---

*Email addresses:* [iakes.goenaga@ehu.eus](mailto:iakes.goenaga@ehu.eus) (Iakes Goenaga),  
[xlahuerta001@ikasle.ehu.eus](mailto:xlahuerta001@ikasle.ehu.eus) (Xabier Lahuerta), [aitziber.atutxa@ehu.eus](mailto:aitziber.atutxa@ehu.eus)  
(Aitziber Atutxa), [koldo.gojenola@ehu.eus](mailto:koldo.gojenola@ehu.eus) (Koldo Gojenola)

vised Machine Learning approaches, using state of the art algorithms like the Perceptron and transfer learning-based Neural Networks.

**Results.** The paper presents a detailed evaluation of the three approaches on two different hospitals. Overall, the best system obtains a 93.03% F-score for section identification. It is worth mentioning that this result is not completely homogeneous over all section types and hospitals, showing that cross-hospital variability in certain sections is bigger than in others.

**Conclusions.** As a main result, this work proves the feasibility of accurate automatic detection and standardization of section blocks in clinical narratives, opening the way to interoperability and secondary use of clinical data.

*Keywords:* Section Identification, Interoperability, Electronic Discharge Summaries, HL7 Clinical Document Architecture

---

## 1 1. Introduction

2     The outstanding advancement of Machine Learning (ML) technologies  
3 (e.g., Deep Learning) enable us to more efficiently harness the large amounts  
4 of data collected through healthcare processes such as clinical narratives in  
5 electronic health records (EHR) as well as electronic discharge summaries  
6 (EDS). EHRs contain a lifetime record of the patient’s complete medical  
7 history, diagnoses and treatment, medications, allergies and immunizations,  
8 as well as radiology images and laboratory results [1]. EDSs are an essential  
9 document to communicate patient journey and care planning regarding *an*  
10 *hospitalization episode* to the next practitioner [2]<sup>1</sup>. In 2016 the proportion  
11 of primary care practices using electronic clinical records was about 80% on

---

<sup>1</sup>Some authors use these two terms interchangeably.

12 average across 15 EU countries [3], and in 2020 in the US the percentage  
13 is of 96% [4]. Digitalization of healthcare systems is contributing to the  
14 improvement of clinical and translational studies, and interoperability and  
15 information exchange between healthcare systems is more necessary than  
16 ever. For that reason, public policies and recommendations are pushing onto  
17 that way [5, 6, 7].

18 There is an increasing interest for integrating heterogeneous health infor-  
19 mation for different reasons: to facilitate the cross-border interoperability of  
20 information among healthcare systems, federal states and countries to ensure  
21 that citizens can securely access and exchange their health data wherever they  
22 are, and also to make digital health information more usable to the bedside  
23 and beyond [5, 6]. Several standards as openEHR [8], HL7-FHIR [9], HL7  
24 CDA/CCR [10] are examples of this standardization effort.

25 However, despite the wealth of information contained in the clinical nar-  
26 ratives, interoperability and extraction of relevant information from docu-  
27 ments remains a challenge. Although the aforementioned standards exist, so  
28 far they have not been widely adopted, and even if so, the healthcare system  
29 at large still has a huge amount of untapped legacy clinical text.

30 Healthcare systems provide guidelines for writing clinical documents,  
31 which for operative reasons typically follow some minimal principles to ensure  
32 the optimal interactions between health professionals and patients like SOAP  
33 (Subjective, Objective, Assessment, Plan), or APIE (Assessment, Plan, Im-  
34 plementation, and Evaluation) [11, 12]. Some systems assume that these  
35 principles are best reflected by using free text, due to flexibility to express  
36 anything that the health-care providers need to record. On the opposite

37 extreme, some impose structured or semi-structured clinical documents in  
38 sections, where each section is a main block of information. In all cases,  
39 the automated processing of clinical texts is hampered by ambiguity, lexical  
40 variety, use of abbreviations, errors due to mistakes, redundancies, etc.

41 Under this scenario, this work presents a first approach towards auto-  
42 matically standardizing Spanish EDSs following the HL7 Clinical Document  
43 Architecture (CDA) R2 template for Discharge Summaries [10] for both help-  
44 ing interoperability and secondary use of Electronic Discharge Summaries.

45 The HL7 CDA R2 template contains a set of clinically relevant sections,  
46 and part of this standardization task is known as Section Identification. It is  
47 defined in [13] as *detecting the boundaries of text sections and adding seman-  
48 tic annotations*. They define a section as *a text segment that groups together  
49 consecutive clauses, phrases or sentences that share the description of one  
50 dimension of a patient, patient’s interaction or clinical findings*. A section  
51 can be marked explicitly, through structural demarcations (headings or sub-  
52 headings), or it can exist implicitly. The main assumption for making this  
53 identification is that unstructured texts have an explicit or implicit structure.

54 Besides its relevance in terms of standarization and interoperability, sec-  
55 tion identification provides a deeper understanding of EDSs, for instance, by  
56 recognizing the section in which a medical entity is located. The same med-  
57 ical condition found in the “past personal medical history” or in the “family  
58 medical history” section might lead to different conclusions. Several works  
59 on secondary use of EHRs and EDSs have shown that section identification  
60 can be helpful for a variety of tasks [14] such as entity recognition [15], co-  
61 hort retrieval [16] and temporal relation extraction [17], and can help in most

62 automatic medical processing tasks, as ICD-10 coding [18, 19, 20, 21]. This  
 63 issue is rapidly becoming an important topic in both academia and industry.

	9027431      XX-XX-XXXX
<b>H</b>	66 años. VARON. MC: REFERENCIADO EN EL INFORME. INFORME AL ALTA :
	Paciente de 66 años. No alergias medicamentosas conocidas.
<b>MH</b>	A. PERSONALES: Enfermedad de Crohn diagnosticada en 1997 con afectación de íleon terminal (A3L1B2) por cuadros suboclusivos resueltos con enfermedad de íleon terminal asociada a mesenteritis fibrosa. Artrosis dorso-lumbar. Cirugía de hernia inguinal. Ci Tto: Dacortin 5: 1-0-0; Pariet 20: 1-0-0; Pentasa:1-1-1; Kilor 0-1-0, Clinutren: 2/día.
<b>CI</b>	E. ACTUAL: Acude a Urgencias por dolor abdominal generalizado con febrícula, sin tiritona, sin náuseas ni vómitos. Sin alteración del ritmo intestinal. Con pauta descendente de corticoides, después del último ingreso por cuadro suboclusivo.
<b>E</b>	EXPL. FÍSICA: Paciente consciente, orientado, colaborador. Buena coloración de piel y mucosas. Cuello: no adenopatías cervicales. AC: rítmica sin soplos. No roncus ni crepitantes. Abdomen: distendido, timpánico. Peristaltismo ausente. Blumberg negativo. EeII: no edemas maleolares. PPP.
<b>EC</b>	RX ABDOMEN: Sugestivo de suboclusión intestinal. Se objetivan dos asas de delgado con niveles hidroaéreos incluso en la cámara gástrica. ANALÍTICA AL INGRESO. Urea, Creatinina, GPT, Amilasa dentro de límites normales. Leucocitos14.400. Segmentados 80 %. TP 100 %. Plaquetas 470.000. Hb 13.5. Hto 41.7 %. PCR 14.2.
	ANALÍTICA AL ALTA: GOT, GPT, Gamma GT, FA, Bilirrubina total, Amilasa, LDH, Colesterol total, Triglicéridos, Na, K dentro de límites normales. Alfa1 antitripsina 169, Albúmina 2.4. PCR 11.6. Fe 18. Transferrina 241. IS 5.3 %. Ferritina 24. Vit B12 247. A fólico 6.3. Hb 11.1. Hto 34.4 %. VCM 78.8. Plaquetas 367.000. Segmentados 75 %. VCM 7.
<b>EV</b>	EVOLUCIÓN Y PROCEDIMIENTOS: Se trata de paciente con enfermedad de Crohn, con afectación ileal y mesenteritis retractoril que ingresa por cuadro suboclusivo, instaurándose tto. corticoideo siendo dado de alta con disminución progresiva de dicho tratamiento. Estando en tto. con 5 mg de Dacortin, ingresa de nuevo con cuadro suboclusivo. Se indica colocación de sonda nasogástrica para aspiración intermitente rechazando el paciente. Se inicia el tratamiento corticoideo iv a dosis plenas, mejorando la clínica del paciente.
<b>D</b>	DIAGNÓSTICO: - CUADRO SUBOCLUSIVO INTESTINAL POR ENFERMEDAD DE CROHN CON AFECTACIÓN ILEAL (A3L1B2) Y MESENERITIS RETRACTI
<b>T</b>	TRATAMIENTO: - Dacortin 60: 1-0-0 durante 1 semana bajando 10 mg cada 10 días hasta 10 mg que mantendrá 15 días más y luego 5 mg 15 días más y suspender.

Figure 1: Example EDS and its sections (H: Heading, MH: Medical History, CI: Current Illness, E: Exploration, EC: Complementary Exploration, EV: Evolution, D: Diagnosis, T: Treatment).

64        Given the difficulty in accurately extracting data from text, most non-  
65 research use of EHR and EDS data rely only on structured data. However,  
66 clinical notes contain highly valuable information not found in strictly struc-  
67 tured fields and, moreover, they give access to volumes of data that are  
68 orders of magnitude bigger and, consequently, improving retrieval accuracy  
69 from text would have great value.

70        In this paper, we will explore the task of section annotation in EDSs  
71 written in Spanish (see Figure 1). We will experiment with three different  
72 approaches, ranging from a knowledge-based solution by means of manually  
73 constructed rules to supervised Machine Learning approaches, including the  
74 structured Perceptron algorithm and Deep Neural Networks. The paper will  
75 present a detailed evaluation of the three approaches and, as a main result,  
76 will prove the feasibility of automatically detecting section blocks in EDSs.

77        The main contributions of this work are:

- 78        • We describe an annotation format for EDSs that defines the section  
79        structure of a document. We have evaluated its feasibility annotating  
80        a dataset comprised of 300 documents and have measured a high inter-  
81        annotator agreement.
- 82        • We implement three different approaches to automatic section identifi-  
83        cation, including a rule-based method, the Perceptron online learning  
84        algorithm and Neural Networks.
- 85        • We conduct exhaustive experiments to explore the contribution of each  
86        method, also giving a detailed analysis of the strengths and weaknesses  
87        of the proposed approaches.

88 The remainder of this paper is structured as follows. Section 2 discusses  
89 related work. The resources and corpus are presented in Section 3. Section 4  
90 sketches the main results, while Section 5 provides an analysis of the results  
91 including a comparison of the different approaches as well as an estimation  
92 of the system’s ability to generalize across hospital settings and a qualitative  
93 evaluation of the encountered errors. Finally, Section 6 summarizes the main  
94 conclusions and future work.

## 95 **2. Background**

96 Pomares-Quimbaya et al. [13] reviewed several studies on clinical section  
97 identification, which varied on the kind of narrative, the type of section, and  
98 the application. The paper examines the characteristics of systems using a  
99 strategy for section identification, the methods used to identify implicit or  
100 explicit sections with different degrees of success, and the main application  
101 scenarios and contexts that have been used with good performance. From the  
102 technical point of view, the methods were classified into rule-based methods  
103 (59%), machine learning methods (22%) and a combination of both (19%).  
104 According to the authors, hybrid methods showed the best performance. 46%  
105 of the studies were able to identify explicit (using headings) and implicit  
106 sections. Regarding the language of application, most of the works (78%)  
107 were intended for English texts.

108 Arnold et al. [22] present SECTOR, a model to segment documents into  
109 sections, under the hypothesis that topics, learned in an unsupervised way,  
110 characterize semantically coherent text segments (sections). Their deep neu-  
111 ral network architecture learns a latent topic embedding over a document, in

112 order to classify local topics and to segment a document at topic shifts. They  
113 report a 56.7% F-score for segmentation and classification in the domain of  
114 diseases. Although the approach seems promising, its main inconvenient for  
115 our task is that, as topics are learned in an unsupervised manner, the topic  
116 clusters do not fit well with the nine HL7 section types of our documents,  
117 because topic clusters can be either finer or more coarse-grained.

118 Choi et al. [23] claim that the structure underlying EHR data improves  
119 the performance of prediction tasks such as heart failure prediction. As most  
120 EHR data do not always contain complete structure information or is com-  
121 pletely unavailable, they experiment alternatives to the baseline consisting  
122 of treating EHR data as a flat-structured bag-of-features. The proposed  
123 model outperformed the baseline approach for various prediction tasks such  
124 as readmission and mortality prediction, indicating that the detection of EHR  
125 structure is beneficial for many tasks.

126 Rosenthal et al. [24] developed a system to detect sections in EHRs, based  
127 on different architectures: an RNN based system and a transfer based system  
128 using BERT. To overcome the lack of annotated data they propose to use  
129 for training purposes sections learned from medical literature (journals, text-  
130 books, web content). They conclude that out of domain clinical literature  
131 is helpful when there is not enough EHR data, but its contribution is not  
132 significant with bigger sizes of the in-domain annotated dataset. Their system  
133 did not exploit the structure of the document, that is, the fact that sometimes  
134 sections are ordered in a canonical order (i.e., first the *Chief complaint*, then  
135 the *Antecedents*, ...), which we plan to use in our approach, as it can be  
136 helpful in deciding section types.



137 Rush et al. [25] solve the section identification problem using a CRF  
138 classifier to mark each token as belonging to a section header, and then  
139 they apply a rule-based post-processing module to structure the annotated  
140 sections. Comparing to our work, they do not perform normalization and  
141 therefore the number of sections they identify is not fixed. In their system,  
142 similar section headers are considered different, while our aim is to normalize  
143 each section into a set of nine HL7 section types.

144 Apart from the medical domain, other areas like legal decision-support  
145 systems also leverage the content structure of documents. For example,  
146 Branting et al. [26] exploit structural and semantic regularities in law case  
147 corpora to identify textual patterns that have both predictable relationships  
148 to case decisions and explanatory value for legal decision support and ex-  
149 plainable outcome prediction.

150 To summarize, we can see that the identification of sections is currently a  
151 promising area of active research, specially for languages other than English.  
152 Historically, rule-based methods have been the most widely used approach,  
153 although the recent emergence of new ML and Deep Learning techniques that  
154 have revolutionized the state of the art on many tasks also presents avenues  
155 for new developments.

### 156 **3. Materials and Methods**

157 In this section we will explore all the corpora and tools we have used  
158 in order to carry out the experiments. In the first part (section 3.1), we  
159 present the annotated corpus, the defined annotation model and the inter-  
160 annotator agreement. Section 3.2 gives a description of the large unannotated

161 corpora used as an additional resource to derive a language model for the  
162 Deep Learning approach (see subsection 3.3.3). Then, section 3.3 describes  
163 the three approaches that have been followed for the automatic annotation  
164 of sections in EDSs.

### 165 *3.1. Annotated Corpus*

#### 166 *3.1.1. The EDS Corpus*

167 We chose to analyze clinical reports of long-term hospital discharges from  
168 two hospitals of the Osakidetza Health System, the Galdakao-Usansolo and  
169 Basurto hospitals. Discharge documents are issued by the responsible doctor  
170 in a health center at the end of each patient’s healthcare process, specifying  
171 the patient’s data, a summary of their clinical history, the healthcare activity  
172 provided, diagnosis and therapeutic recommendations. They are documents  
173 of great importance within the clinical history, containing the summary of  
174 the care provided to the patient during the hospitalization episode. The  
175 recipients are different users with diverse interests, including the patient,  
176 his/her family, the primary care physician and the specialist physician.

177 A set of 300 documents was selected for manual annotation (see Table  
178 3) divided evenly between the two hospitals. As each EDS typically can  
179 contain most of the nine section types that we have defined, this corpus,  
180 albeit of a moderate size, can give a sufficient amount of data (more than  
181 2,000 instances of the different section types) for training and evaluation.

#### 182 *3.1.2. Definition of the Main Section Types*

183 As mentioned in the introduction, we have followed the HL7 CDA R2  
184 recommendations proposed for EDSs. This standard requires EDSs to be

185 minimally structured in at least three main sections: Hospital Course Sec-  
186 tion, Discharge Diagnosis Section and Plan of Treatment Section. Never-  
187 theless, there are 22 other sections that are optional and try to cover the  
188 heterogeneity of information captured in EDSs. This optionality allows each  
189 healthcare system to select those sections that better accommodate to the  
190 reported information. In our case and following [27] we selected 9 sections  
191 out of those 22. Table 1 summarizes the adopted HL7 CDA R2 sections spe-  
192 cially recommended for Electronic Discharge Summaries along with a short  
193 description and the nomenclature we will employ all over the paper.

194 Although ideally all EDSs could contain each and all of the listed sec-  
195 tions, in practice this just does not happen most of the times. It is usual to  
196 find summaries with less elements and it can also happen to find some sec-  
197 tions more than once in a document, possibly when the discharge summary  
198 includes more than one episode from one patient.

199 Regarding the order, although the given description of section types can  
200 be considered as a canonical ordering of the elements in an EDS, there is a  
201 great variability. Except for the heading, that appears almost always in the  
202 first position, the rest of the sections can be found in different parts of the  
203 document. For example, even if it is common to find the diagnosis and treat-  
204 ment at the end of the EDS, some practitioners tend to move them towards  
205 the beginning of the document. Even when many sections are marked by  
206 an explicit heading, there are several challenges related to the detection of  
207 section boundaries:

- 208 • Variability of section headings. Although the standard definition could  
209 suggest that all the headings are naturally defined by a common term,

HL7 CDA R2 name	Abbreviated name	Description
-	Header (H)	Not a HL7 CDA R2 section but included to capture the header, which can be highly specific to each hospital
Chief complaint	Chief complaint (CC)	This section is similar a press headline, and it briefly contains the answers to who, what, where, why and when
Past Medical History	Medical History (MH)	Past symptoms, medications, diseases or procedures. Sometimes there are specific subsections for <i>Family history</i> or <i>Personal history</i> .
History of Present Illness	Current Illness (CI)	A detailed description of the issues presented in the chief complaint section.
Vital Sign Section	Exploration (E)	This section describes observations, including the vital signs, muscle power and examination of different organs, especially ones that might be related to the symptoms.
Hospital Discharge Studies Summary	Compl. Exploration (EC)	Additional, specialized tests, like ECG, or a radiography.
Review of Systems	Evolution (EV)	Evolution of the patient during the hospitalization.
Discharge Diagnosis	Diagnoses (DI)	Main and secondary diseases diagnosed by a medical practitioner.
Plan of Treatment	Treatment (T)	Medications, procedures and recommendations for this patient

Table 1: Brief description of the HL7 CDA R2 adopted sections.

210 there is a great variability, corresponding to the use of synonyms, ab-  
211 breviations and variations. Additionally, in some cases section headings  
212 can be misleading or ambiguous, and the content of the text accompa-  
213 nying the heading must be taken into account in order to disambiguate  
214 the text.

215 • Implicit sections. Apart from the headings, sections can also be identified by looking to the body text. For example, the description of  
 216 measures like Sodium or Potassium will typically appear at the Exploration section. This information is also useful to detect sections, and  
 217 is the only possibility when the section heading is not present. Table 2 shows how a considerable proportion of sections do not have an explicit  
 218 section header. For example, only a 30% and 41% of the *Chief Complaint* and *Complementary exploration* sections, respectively, are  
 219 explicitly marked. Some others, although they have an explicit heading most of the times, show a great variability (e.g., *EXPLORACION* in  
 220 Table 2).

```

  <?xml version = "1.0"?>
  <sections>
  <section id="1" type="ENCABEZADO" str="" offset="1-6"></section>
  <section id="2" type="ANTECEDENTES" str="A. PERSONALES" offset="7-10"></section>
  <section id="3" type="ENFERMEDAD ACTUAL" str="E. ACTUAL" offset="11-13"></section>
  <section id="4" type="EXPLORACION" str="EXPL. FÍSICA" offset="14-19"></section>
  <section id="5" type="EXPLORACION COMPLEMENTARIA" str="RX ABDOMEN" offset="20-29"></section>
  <section id="6" type="DIAGNOSTICO" str="IMPRESIÓN DIAGNÓSTICA" offset="30-32"></section>
  <section id="7" type="EVOLUCION" str="EVOLUCIÓN Y PROCEDIMIENTOS" offset="33-36"></section>
  <section id="8" type="DIAGNOSTICO" str="DIAGNÓSTICO" offset="37-39"></section>
  <section id="9" type="TRATAMIENTO" str="TRATAMIENTO" offset="40-50"></section>
  </sections>
  
```

Figure 2: Example of EDS annotation corresponding to Figure 1.

226 We chose to use a stand-off annotation based on XML. For example,  
 227 Figure 2 presents the annotation document for the EDS presented in Figure  
 228 1. Each section is described by an XML element containing attributes for  
 229 the section type, the string (if any) that indicates the start of the section

Section	Examples
<b>ENCABEZADO</b> (HEADING)	9027431 16-04-09 66 años.
	VARON. MC: REFERENCIADO EN EL INFORME.
<b>MOTIVO DE CONSULTA</b> (CHIEF COMPLAINT)	MOTIVO DE INGRESO ...
	MOTIVO DE CONSULTA ...
	Paciente que ingresa procedente de ...
	Paciente varón de 47 años que ingresa para ...
	Varón de 87 años ingresado desde ...
	Varon de 63 años que consulta por ...
	MI ...
<b>ANTECEDENTES</b> (MEDICAL HISTORY)	ANTECEDENTES PERSONALES ...
	A.PERSONALES ...
	AP ...
	A. Personales ...
	A.P. ...
	A. PERSONALES ...
	Paciente de 65 años de edad con antecedentes ...
<b>EXPLORACION</b> (EXPLORATION   PHYSICAL EXAMINATION)	EXPLORACION GENERAL ...
	EXPLORACION ...
	Exploración física ...
	EXPLORACION VASCULAR ...
	EXPLORACIÓN ORL ...
	EXPLORACIÓN PSICOPATOLÓGICA ...
	EXPLORACIÓN CLÍNICA EN LA UNIDAD ...
<b>EXPLORACION COMPLEMENTARIA</b> (COMPLEMENTARY EXPLORATION)	PRUEBAS COMPLEMENTARIAS ...
	EXPLORACION COMPLEMENTARIA ...
	ECG ...
	RX ABDOMEN ...
	ANALITICA AL ALTA ...
	ANALÍTICA EN URGENCIAS ...

Table 2: Examples of some instances of the beginning of identified Section Types.

230 and the offset (in lines) corresponding to the text of the section<sup>2</sup>.

<sup>2</sup>Although the string attribute was useful for annotators when discussing any annota-

231 After the main section types to be annotated were defined, a set of docu-  
232 ments from two different hospitals was annotated by two annotators. Table  
233 3 describes the three-way data split (training, development and test). In  
234 order to minimize the annotation effort, the number of annotated documents  
235 could not be too large but it should also provide enough data for training  
236 and evaluation. Taking these considerations into account, a corpus of 300  
237 documents was randomly selected. Regarding the split of the dataset into  
238 three subsets corresponding to training, development (or validation), and the  
239 final test, the validation and test sets should contain enough instances of each  
240 section type for the evaluation to be significant. For this reason we decided  
241 that each of the three subsets would contain 100 documents, different from  
242 classical data splits (e.g., 70%, 15% and 15% for training, development and  
243 test, respectively). Figure 3 shows the distribution of sections in both the  
244 train and development splits. Note that all sections are represented in both  
245 splits and their distribution is similar.

---

tion disagreement, in fact the attributes that should be obtained by an automatic system will be the section type and its location (line offset) in the document.

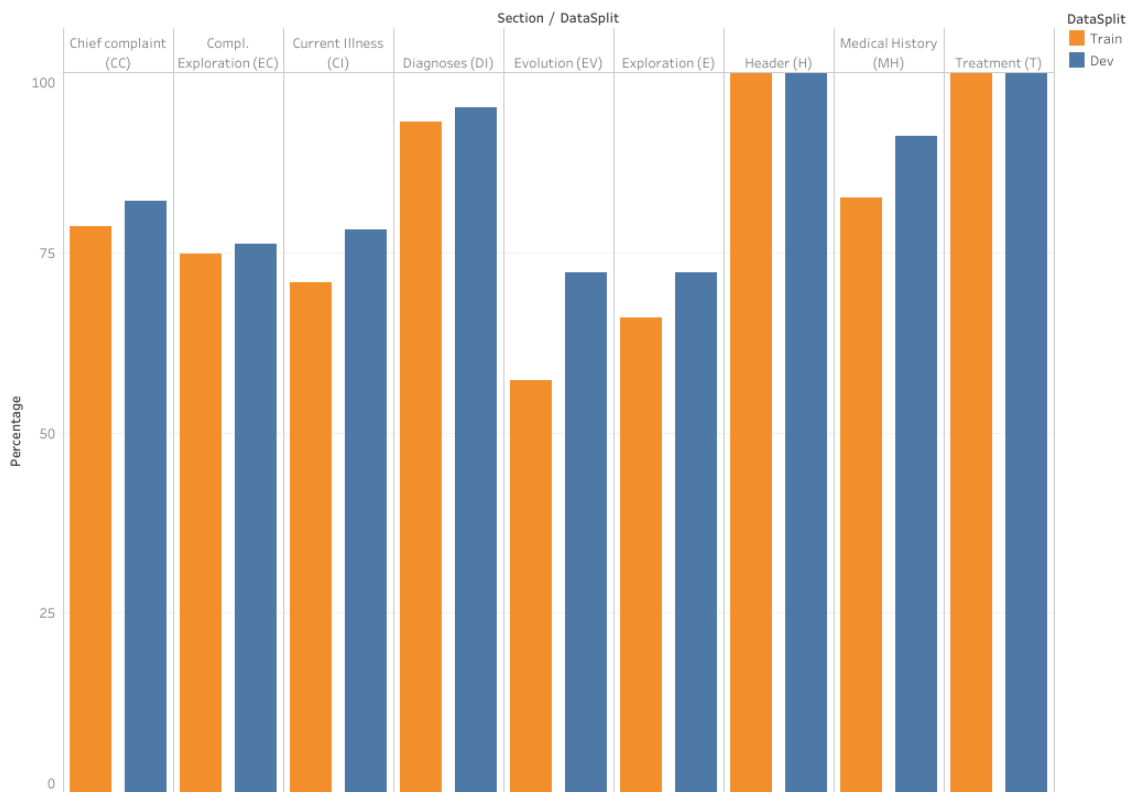


Figure 3: Distribution of sections in both train and development splits. Section chief complaint, for instance, is present in 82% of the development split, while slightly less in the training split (79%).

	<b>Documents</b>	<b>Sections</b>	<b>Tokens</b>
<b>training</b>	100	744	47,449
<b>development</b>	100	786	48,461
<b>test</b>	100	754	59,119

Table 3: Details of the annotated corpus.



### 246 3.1.3. *Inter-Annotator Agreement*

247 The annotation process of the corpus was performed by annotation ex-  
248 perts. Annotators followed an iterative process of training until a high inter-  
249 annotator agreement was reached. The final agreement measure was cal-  
250 culated on a set of 25 EDSs that were doubly annotated by two different  
251 annotators, reaching a pairwise agreement of 93.47% Cohen’s Kappa, indi-  
252 cating that the agreement is very high. There were differences with respect  
253 to each section type, ranging from 86% for Diagnosis (lowest agreement) to  
254 100% for some section types, thus reaching a significant agreement for all  
255 section types.

256 Our annotation strategy requires each section type to be matched exactly  
257 while taking into account its content, and additionally returns the first and  
258 last lines of each section. While this strategy might seem an overly stringent  
259 criteria, the task is well defined as evidenced by the high inter-annotator  
260 agreement.

### 261 3.2. *Textual Corpora*

262 Deep learning techniques usually require huge amounts of data. Although  
263 manually annotated data gives the best results, it is very expensive and time  
264 consuming. For that reason, the idea of acquiring useful information in an  
265 unsupervised manner is very attractive, and efficient and effective methods  
266 have been developed. Vectorial representations of words, also known as *word*  
267 *embeddings* [28, 29], that are learned from textual corpora, have proven useful  
268 as an information source for many Natural Language Processing tasks, such  
269 as Part-Of-Speech (POS) tagging, Named Entity Recognition or Machine  
270 Translation, due to their ability to acquire relevant generalizations. These

271 embeddings are learned through solving an appropriate optimization objec-  
272 tive [28] under the assumption that similar words occur in similar contexts.  
273 As a result, vectors of similar words derived from such optimization tend  
274 to reside in the neighborhood in the vector space. There are two kinds of  
275 embeddings, static and contextual. Static embeddings capture in a vectorial  
276 representation information of a word form, while contextual embeddings are  
277 sensitive to context, representing both a word and its context.

278 This way, an unsupervised system can utilize the information based on  
279 word similarity in a manner that associates unseen words with those already  
280 occurring in the annotated corpus, thereby allowing us to cover unseen and  
281 misspelled terms. For instance, *infarct* and *stroke* are similar terms but one  
282 of them may not be in the annotated data set. The resulting word vectors  
283 will be fed to the neural network as input during training (see Figure 4), thus  
284 providing a model of the language that can help obtain better generalizations  
285 and, consequently, increase the recall of the final tool.

286 For this work we have employed heterogeneous embedding information  
287 both static and contextual in order to make the system sensitive to different  
288 granularity and domain specificity. Regarding the granularity, during the  
289 section identification training, adding a character embedding layer allows  
290 the system to learn at the character level. Besides the character embed-  
291 dings learned during the training, we incorporated pre-trained, character-  
292 based embeddings based on fastText [30] trained over the Spanish version  
293 of Wikipedia. Character-based embeddings are able to generalize over n-  
294 grams, enabling the system to take into account prefixes and suffixes as well  
295 as to capture information about the different n-gram variations on the sec-

296 tion heading words. They also generalize over zero-shot words, words that  
 297 do not appear in the training corpus as their building element are characters  
 298 and not words.

299 Table 4 presents the details of the different word embeddings we have  
 300 used for the task. Static embeddings were obtained applying word2vec [30]  
 301 to Electronic Discharge Summaries (50M words), together with pretrained  
 302 embeddings that had been calculated with Wikipedia2Vec [31], representative  
 303 of general domain. Additionally, we also used contextual string embeddings  
 304 [32] we calculated from Electronic Discharge Summaries and Wikipedia .

<b>Technique</b>	<b>Source text</b>	<b>Embedding type</b>	<b>Details</b>
word2vec	EDS	static	window length = 1, dimensions = 300, algorithm = SkipNgram
Wikipedia2Vec	general domain		window length = 5, dimensions = 300, algorithm = Skipgram
FLAIR	EDSs	contextual	layers=1, hidden size = 2,048, sequence length = 250, mini batch size = 32
	general domain		layers = 1, hidden size = 1,024, sequence length = 250, mini batch size = 100

Table 4: Overview of the different embedding types used in this work (static word embeddings and contextual character embeddings).

### 305 3.3. Approaches to Automatic Section Identification

306 In this section we will explain the different approaches we have tried with  
307 the aim of automatically identifying sections in medical records. First of  
308 all, in subsection 3.3.1 we will specify the setup we used for the rule-based  
309 tool that we have developed. After that, in subsections 3.3.2, and 3.3.3 we  
310 will explore the ML algorithms we have employed, the Perceptron and Deep  
311 Learning, respectively. For both ML approaches, we have approached the  
312 task as a sequential learning process [33, 34], where the text is considered a  
313 sequence of tokens, and each token is associated with one tag indicating its  
314 corresponding section. We have used an IOB (Inside, Outside, Begin) tag  
315 model, where the beginning of each section is marked with a B tag (e.g., *B-*  
316 *DIA* for the token starting a diagnosis), the tokens inside a section are marked  
317 with an I tag (*I-DIA* will mark a token inside a diagnosis section), and using  
318 the O tag for elements that do not belong to any section (see Figure 5). This  
319 way, section identification can be viewed as the detection of extended and  
320 long entities. This approach has been successfully used in similar tasks as  
321 the identification of elementary discourse units (text segments consisting of  
322 one or several sentences) in Discourse processing [35] or topic segmentation  
323 [22]. Figure 4 presents an architecture of the system.

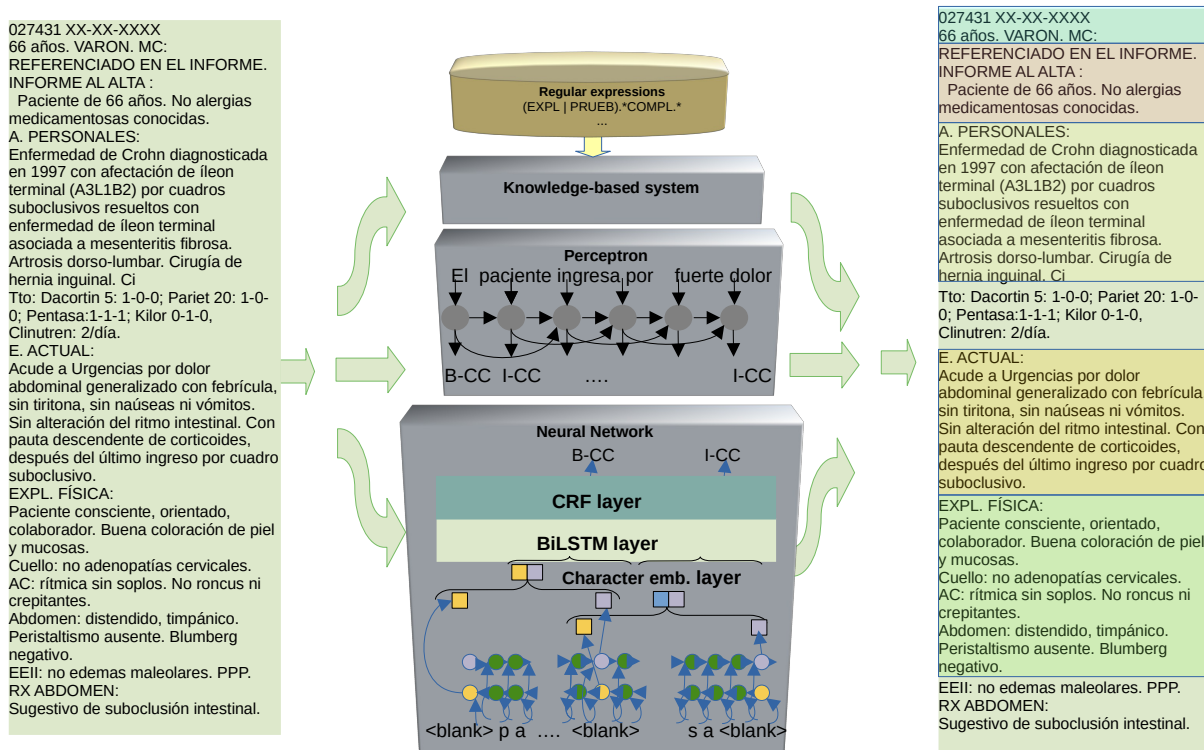


Figure 4: Architecture of the system. Three different approaches have been used: regular expressions, the Perceptron algorithm and neural networks.

### 3.3.1. Rule-Based Approach

Manually defined rules have been used since the early years of Artificial Intelligence, and are still a competitive method to achieve acceptable results. Their downside is the effort needed to include knowledge into the automatic system. Another drawback is their lack of generalization, because a change in the domain may imply a complete re-implementation of the rule system.

Regarding the identification of sections in medical records, this approach

331 has been used in many systems [13, 36], where acceptable results have been  
 332 reported, although in several cases the approach has not been general, but  
 333 rather limited to a reduced set of very specific sections or portions of text.

Table 2 presents several examples of the beginning of different section types. The table shows how there is a high variability difficult to capture using rules, specially with implicit sections with no standard title, like in the *Chief Complaint* and *Complementary Exploration*. Examples (1) and (2) present two rules that try to capture the start of the *Chief Complaint* and the *Current Illness* sections, where the parentheses enclose optional elements. The objective was to cover the different options found in the training set.

- 334 (1) MOTIVO(S) (DE(L)(A)) *INGRESO*|*PETICION*|  
 335 *EXPLORACION*|*ESTUDIO*|*CONSULTA* (*ACTUAL*)
- 336 (2) (*E.*|*ENFERMEDAD*|*SITUACIÓN*|*EPISODIO*|*ESTADO*)  
 337 (*A.*|*ACTUAL*) | *SINTOMATOLOGÍA*

### 338 3.3.2. Machine Learning: Perceptron

339 For the application of ML to section identification, we modeled the prob-  
 340 lem as a sequence to sequence problem. The task consists in learning to  
 341 map from input word sequences  $w_1 \dots w_m \mid w_i \in W$  to output tag sequences  
 342  $t_1 \dots t_m \mid t_i \in T$ .

343 Although some approaches to section identification used sentence se-  
 344 quences as input units [13], we preferred to model this problem using word  
 345 sequences as input units to capture the fact that individual words in the

346 right context are good signals for sections and also to reduce sparsity, be-  
347 cause sentence sequences are more sparse than word sequences. The problem  
348 is cast as the assignment of the correct tag to each token. Although the tag  
349 assignment is made token by token, the final evaluation will be done on the  
350 detection of complete sections.

351 To do so, we employed the Averaged Structured Perceptron algorithm  
352 [37, 38] which combines the Perceptron algorithm for learning linear classi-  
353 fiers with an inference algorithm and converts a classification problem into a  
354 ranking problem. The objective of the algorithm is to find, for each sentence,  
355 the sequence of tags with the maximum score. This prediction decision pro-  
356 cess is divided into a sequence of smaller decisions made from left-to-right.  
357 Thus, at each step there is a word and its context, called *the history*, in  
358 which the local tagging decision is made, namely to predict the tag given the  
359 history. The history can be represented in several ways, using the prefixes of  
360 a given number of previous words, and/or the suffixes, or any other features  
361 that could be relevant for the task and then converted into a feature vector  
362 where each feature will get a weight through the learning process.

363 Formally, the problem can be stated as follows. Given:

- 364 • A sequence of input words  $w_1 \dots w_m$ , for simplicity referred as  $w$ .
- 365 • The sequence of tags  $t_1 \dots t_m$  as  $t$  (this way, the set of possible tags is  
366  $T$ ).
- 367 • In our case the context in which a tagging decision is made is repre-  
368 sented by the history tuple  $h: \langle t_{-2}, t_{-1}, w_{-2}, w_{-1}, w_0, w_{+1}, sx_0, px_0, cap,$   
369  $num, i \rangle$ , where  $t_{-2}$  and  $t_{-1}$  are the previous two tags,  $w_{-2}$  and  $w_{-1}$  are

370 the previous two words at a given position  $i$  (this way,  $H$  corresponds  
 371 to the set of all possible histories).

372  $sx$  and  $px$  correspond to different sizes of word suffixes, (in this work,  
 373  $x$  varying from 2 to 4) and prefixes of  $w_0$ .

374  $cap$  and  $num$  correspond to two binary features to account for capital-  
 375 ization and number status at the current word.

376 The feature mapping function  $\Phi : H \times T \rightarrow \mathbb{R}^d$  maps a history-tag  
 377 pair to a d-dimensional feature vector we mentioned before. The Structured  
 378 Perceptron models  $P(t|w)$  as  $P(t|h; \alpha)$  where  $\alpha \in \mathbb{R}^d$  is a parameter vector  
 379 representing the weight of each feature of  $\Phi$ .  $P(t|h; \alpha)$  is calculated as  $\alpha \cdot$   
 380  $\Phi(h, t)$  and the objective function is:

$$381 \quad \hat{t} = \operatorname{argmax}_t \sum_1^d \alpha_i \cdot \Phi_i(h, t)$$

382 Usually the Viterbi algorithm is applied when used on sequence data,  
 383 in order to efficiently calculate the best tag sequence using dynamic pro-  
 384 gramming. The algorithm is competitive to other options such as maximum-  
 385 entropy taggers or CRFs [33].

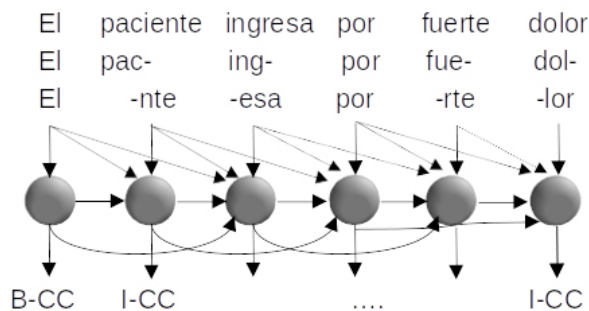


Figure 5: Simplified Architecture of the Structured Perceptron (the upper three rows exemplify the use of word features (first row), 3 letter prefixes and 3 letter suffixes (second and third rows)).



386 We employed our own implementation of this tagger following [37]. We  
387 trained 100 iterations and selected the model corresponding to the iteration  
388 that achieved the best score on the development set. Although the algorithm  
389 achieves a competitive performance compared to state of the art methods,  
390 this approach requires a feature engineering effort to identify, select and  
391 properly encode relevant features.

### 392 *3.3.3. Machine Learning: Neural Networks*

393 In addition to a traditional neural network like Perceptron, we explored  
394 transfer learning methods. In this case, we used FLAIR [39], a bi-directionally  
395 trained Language Model (LM) using Recurrent Neural Networks (RNN),  
396 where the basic element is the character and not the word. Based on its char-  
397 acters, FLAIR generates pre-trained contextual embeddings for each word by  
398 concatenating the hidden state for the last character of the word in the for-  
399 ward neural network and the first character of the word in the backward  
400 neural network, as shown in Figure 6. As described in [39], formally, the ob-  
401 jective function of a character-based LM is to maximize the sum of the logs  
402 of  $P(x_t|x_0, \dots, x_{t-1})$ , that is to say, an estimate of the predictive distribution  
403 over the next character given past characters. FLAIR allows us to com-  
404 bine different types of embeddings by concatenating each embedding vector  
405 to form the final word vector. We employed a combination of embeddings  
406 as previously reported in section 3.2. One of the main advantages of these  
407 methods is that there is no need for feature engineering.

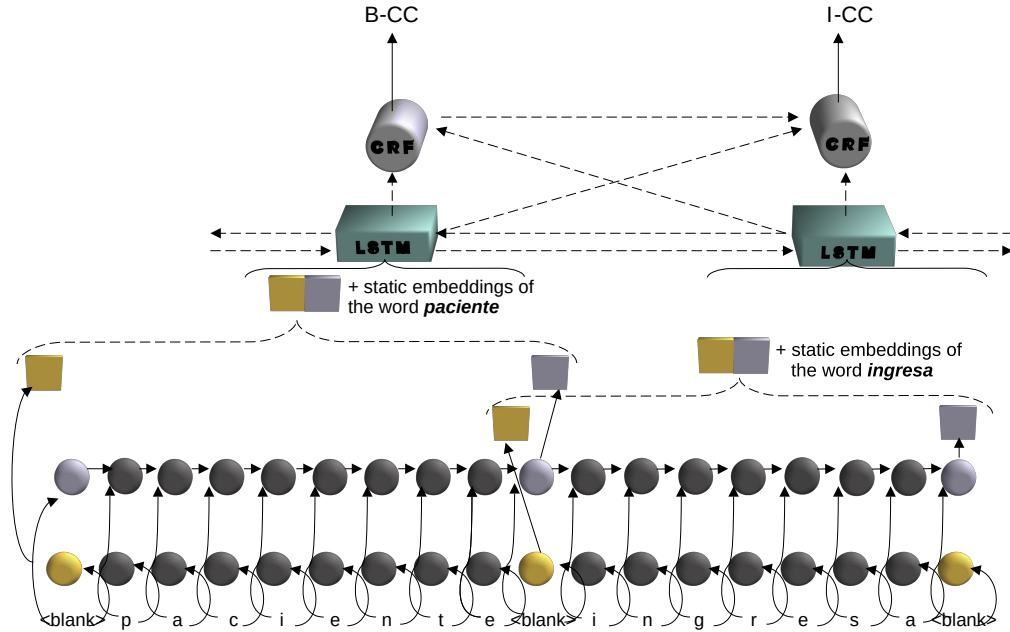


Figure 6: Simplified Architecture of FLAIR (B-CC: start of Chief Complaint, I-CC: continuation of Chief Complaint).

	Micro average			Macro average (per document)		
	Precision	Recall	F-score	Precision	Recall	F-score
<b>Rule-based</b>	52.86	51.63	52.24 (14.8)	58.02	57.24	57.62
<b>Perceptron</b>	87.28	85.18	86.22 (3.1)	87.34	86.10	86.72
<b>Neural Networks</b>	93.40	92.55	93.03 (1.8)	91.77	91.26	91.52

Table 5: Results of the different approaches on Section Identification (margin of error with 95% confidence in parentheses).

#### 408 4. Results

409 For evaluation, employed the standard measures of precision, recall and  
410 F-score defined in formulae (1), where  $TPS =$  *correctly identified sections*,  
411  $FPS =$  *incorrectly identified sections (marked by the automatic tool and not*  
412 *present in the annotated gold standard)* and  $FNS =$  *false negatives, i. e.,*  
413 *present in the gold standard and not detected by the automatic tool*. Table  
414 5 shows the main results, given as micro average (over all instances and all  
415 documents) and macro average (calculating the mean over the scores on each  
416 document).

$$\begin{aligned} Precision &= \frac{TPS}{TPS + FPS} \\ Recall &= \frac{TPS}{TPS + FNS} \\ F\text{-score} &= \frac{2 * Precision * Recall}{Precision + Recall} \end{aligned} \tag{1}$$

417 In predicting the section type, the evaluation has been strict in the sense  
418 that an error has been counted whenever an automatically detected section  
419 did not exactly match with the gold standard section. This was done even  
420 when in some cases there is a high degree of overlapping with a correct  
421 section (e.g., when the system correctly marks a paragraph as belonging to  
422 the Medical History, but it also misses a part of the gold standard section).

423 Looking at Table 5, we see that the rule-based approach gives the lowest  
424 performance (52.24 and 57.62 for micro and macro average, respectively), far  
425 from the Machine Learning approaches, and contrary to our first intuition.  
426 In general, rule-based solutions tend to have a better precision at the cost of a  
427 lower recall, although in this experiment there is not a significant difference  
428 between precision and recall. The example rule (1) obtained the best F-

429 score (71) for the *Chief Complaint* section type, and our successive efforts to  
430 improve it were not successful, because our attempts to boost recall worsened  
431 the precision. Regarding the reason why the rule-based approach gave the  
432 worst results, in Information Extraction usually designing more general rules  
433 gives an increase in recall, while more specific rules tend to improve precision  
434 at the cost of diminishing recall. However, in our particular problem this is  
435 not the case, because the objective of the rules is to exactly match entire  
436 sections. As a consequence, more general rules have a negative effect on  
437 both precision and recall, since incorrectly marking a section produces a  
438 cascade effect in the surrounding sections. This was the cause why, after the  
439 first successful attempts, dedicating more effort to the rules deteriorated the  
440 performance. Thus, we concluded that for this experiment even customized  
441 and carefully designed rules with a time-consuming implementation were not  
442 able to increase precision. Both ML approaches surpass the performance of  
443 the rule-based system, being the neural network based system the best one  
444 by a significant margin. The Perceptron-based system outperforms the rule-  
445 based one by around 30 absolute points, while neural networks give the best  
446 result with 93.03 and 91.52 F-score for micro and macro average, respectively.

447 Figure 7 presents a detailed comparison of the performance of each ap-  
448 proach on the different section types. The rule-based approach presents the  
449 lower results for all section types, although the F-score is high in the *Heading*  
450 section, which can be considered the easiest to detect. Specially bad is the  
451 result for *Exploration*, *Complementary Exploration* and *Evolution*, possibly  
452 due to the fact that these sections are not usually marked by an explicit  
453 heading. Secondly, the Perceptron based system ranks after the neural net-

454 work based one, but their F-score is similar for the *Diagnosis* and *Treatment*  
455 section types, which can be considered the most important ones from the  
456 point of view of the automatic processing of EDSs.

457 Overall, we can also see how some sections are harder to detect than  
458 others. This fact can drastically affect the rule-based system, with big dif-  
459 ferences according to each section type, and it is related to the difficulty of  
460 finding patterns for sections where the headings are absent or also with sec-  
461 tions where the headings present a high variability. The differences, albeit  
462 smaller, also appear with the Perceptron, which although it is an automatic  
463 Machine Learning algorithm, requires an explicit definition of features based  
464 on words, suffixes, prefixes or capitalization (feature engineering). In this re-  
465 spect, the diagnoses and treatment sections present the best results, possibly  
466 due to the fact that their headings are more predictable. Finally, the neu-  
467 ral network system is able to detect the sections without an explicit feature  
468 definition. As this system is based on left-to-right and right-to-left vector  
469 encodings of the processed text, these systems are able to learn not only  
470 from the headings, but also from the vocabulary inside the sections. This is  
471 the reason why this system outperforms the other two in the sections with  
472 the lowest proportion of explicit headings, like *Exploration* and *Complemen-*  
473 *tary exploration*, where the elements appearing inside the section text, like  
474 procedures (ECG, X rays, ...) or clinical measures (sodium, potassium, ...)  
475 can help to decide which section is being examined.

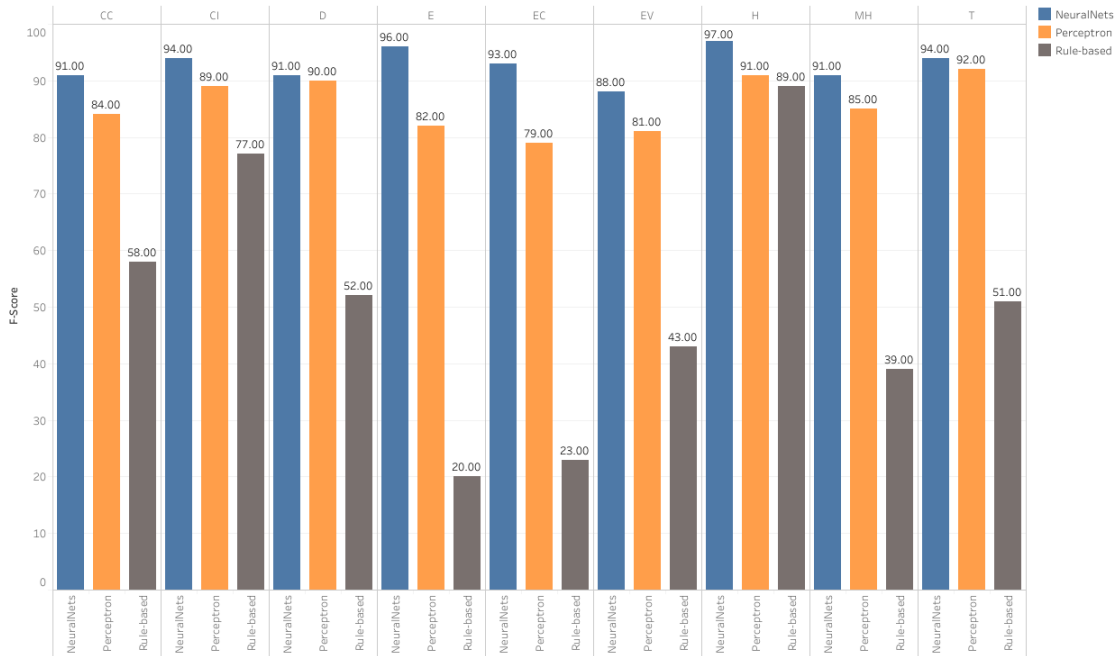


Figure 7: Comparison (F-score) of the three approaches on each section type (H: Heading, CC: Chief Complaint, MH: Medical History, CI: Current Illness, E: Exploration, EC: Complementary Exploration, EV: Evolution, D: Diagnosis, T: Treatment).

476 Examining the results for the best (neural) system on individual sections,  
 477 we can see that EV(olution) presents the lowest result (88% F-score), other  
 478 sections like C(hief) C(omplaint), M(edical) H(istory) and D(iagnosis) give  
 479 better results (91% F-score), and the remaining sections present higher ac-  
 480 curacies.

## 481 5. Discussion

482 In the next subsections we will first (subsection 5.1) look at a new set of  
 483 experiments to address the effect of applying our system to a new hospital,

484 as usually the writing of EDSs can vary greatly from one hospital to another  
485 belonging to the same Health System. Finally, subsection 5.2 presents an  
486 analysis of the main error sources.

### 487 *5.1. Cross hospital generalization*

488 Usually, ML systems tend to obtain good results when the domain of  
489 application is the same as the one used for training but, when moving to a  
490 different scenario or domain, the results can degrade drastically. We wanted  
491 to test the effect of changing the environment of application and, knowing  
492 that many times writing styles can vary from one hospital to another, we  
493 measured the effect of training using data from one hospital and testing  
494 on EDSs from a different hospital. In our case, our data came from two  
495 different hospitals from the same hospital system (the Basque Health System,  
496 Osakidetza). Figure 8 shows the results when testing on data from a hospital  
497 when the training data belongs to the same or a different hospital. The  
498 experiments were performed using the best system in Section 4, the neural  
499 network based one. As could be expected, the best results are obtained when  
500 the training and test sets belong to the same hospital (two left-side bars in  
501 each column in Figure 8), and the scores worsen when the system is trained  
502 on data from a hospital and applied to the other one (two right-side bars).

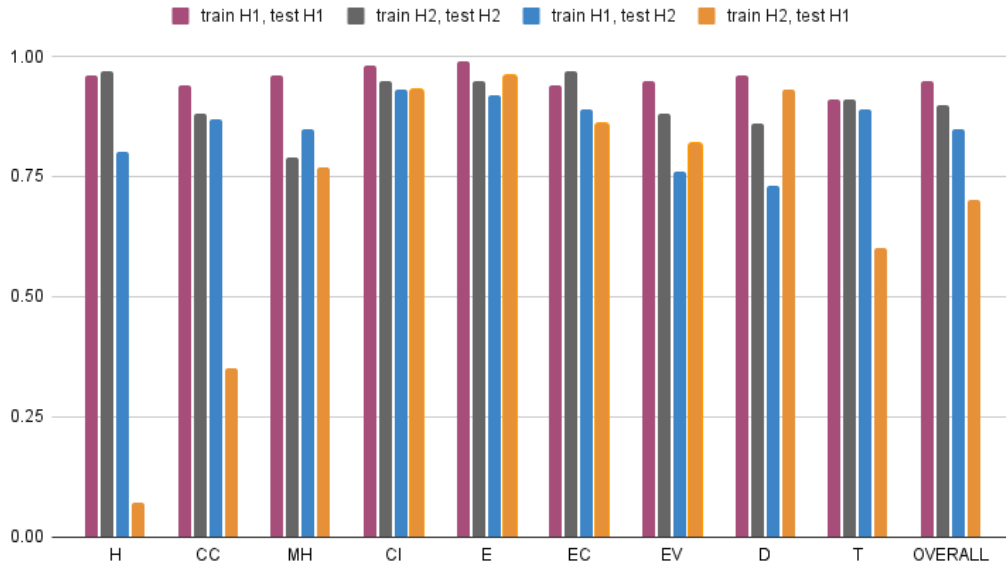


Figure 8: Effect of training and testing on the same or a different hospital (H1: Galdakao-Usansolo hospital, H2: Basurto hospital), measured by F-score.

503 The difference is significant in almost all types of sections. Specially  
 504 relevant is the effect of the system trained on hospital H2 and tested on  
 505 hospital H1 for the Heading section type (H column), where the F-score is  
 506 very low. We examined the results and concluded that this happens mostly  
 507 because the headings show a great variation, added to the fact that the  
 508 data present in the headings is generated automatically most of the times,  
 509 including record numbers or dates, and this implies that they can be different  
 510 enough to confuse an automatic system. Surprisingly, this does not happen  
 511 in the opposite direction, meaning that the data from hospital H1 shows  
 512 more variability and is useful to account for the instance types of hospital  
 513 H2. The difference is also significant for the second section type (Chief



514 Complaint, CC), although less drastic. This was due to a cascade effect  
515 as a result of applying a sequence to sequence approach, as the errors in  
516 delimiting the first section of the document frequently are carried from one  
517 section to the next one. Finally, for some section types, like E(xploration),  
518 EV(olution) and D(iagnostic), we can see how applying a system trained on  
519 a different hospital can outperform the system based on data from the same  
520 hospital. This can be due to the fact that one hospital agrees more with the  
521 conventions of the other hospital for these section types.

## 522 *5.2. Error Analysis*

523 We looked at the errors given by the different systems. For simplicity, we  
524 will only examine the results of the best system based on neural networks.  
525 An examination of the divergences between the output of the system and the  
526 gold standard showed us the main causes of error:

- 527 • Errors given by the inherent difficulty of spontaneously written section  
528 headings. Although explicit headings are an important clue to delimit  
529 sections, the variability of their writings together with the limited size  
530 of the training set (100 documents, which means that there are at most  
531 100 instances of each section type) is a source of errors.
- 532 • Implicit sections. Some types of sections have a majority of instances  
533 without an explicit section heading, which means that the section must  
534 be detected using its content words (see Table 2).
- 535 • Mixed sections. Although the annotators have decided the exact scope  
536 of each section with a high agreement, the use of unstructured and

537 spontaneously written EDSs gives the writers freedom to describe any  
538 concept in different places. As an example, the section corresponding  
539 to the *Medical History* can contain passages related to past diagnoses,  
540 treatments and explorations, which can pose a challenge for an auto-  
541 matic system.

542 • A special case of mixed sections can be the confusion between two  
543 related section types:

544 – *Chief Complaint* and *Current Illness*. These two sections present  
545 the most diffuse definition [27], and are the cause of several errors.

546 – *Exploration* and *Complementary Exploration*. Although the defi-  
547 nition of each section is precise, sometimes physicians mix them  
548 in the same block or paragraph.

549 In Section 3.1.2, we mentioned that the ordering of section types shows  
550 a great variability. In order to measure its impact on the results, we split  
551 the test set in two subsets. The first subset corresponds to the documents  
552 that follow the canonical order (26% of the documents), while the rest of the  
553 documents conform the second subset (non-canonical order and/or missing  
554 sections, 74% of the documents). Since our sequence learning-based methods  
555 depend on the ordering for predicting the next token, this has an effect in the  
556 IOB-labeling prediction, with a F-score of 95.40 for the canonical documents  
557 and 89.81 for the non-canonical ones.

558 Figure 9 presents the main types of mistakes made by the automatic  
559 tool. It shows how the errors are concentrated in some sections, like Chief  
560 Complaint (CC), Medical history (MH) and Diagnosis (D). Overall, the dis-

561 tinction of different sections is reflected in the text by means of different clues,  
 562 ranging from semantics (the content of each section) to syntax (e.g., use of  
 563 section headings and separate paragraphs or text blocks for each section)  
 564 but, in most of the errors, these conventions do not hold, and this causes the  
 565 automatic tool to find an additional difficulty.

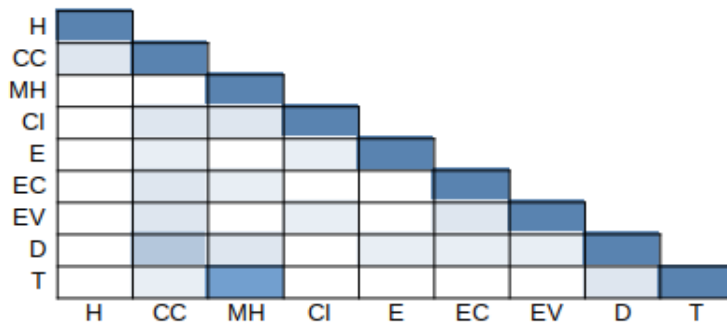


Figure 9: Confusion matrix, where darker green means a higher frequency, of each instance (H: Heading, CC: Chief Complaint, MH: Medical History, CI: Current Illness, E: Exploration, EC: Complementary Exploration, EV: Evolution, D: Diagnosis, T: Treatment).

## 566 6. Conclusion

567 We present a system for Section Identification in Discharge Summaries  
 568 written in Spanish. We have adopted an annotation model based on H7 CDA  
 569 R2 for Electronic Discharge Summaries (EDS) of the Spanish Health System,  
 570 and we have applied it to manually annotate a corpus of 300 EDSs, obtaining  
 571 a high inter-annotator agreement.

572 We have evaluated the contribution of different rule-based and Machine  
 573 Learning approaches and study the strengths and weaknesses of each option.  
 574 Most previous works have used section identification as an auxiliary module

575 for carrying on clinical processing, relying on a rule-based approach. How-  
576 ever, our results show that section identification is a task on its own, where  
577 simple methods do not obtain the best results. The Machine Learning sys-  
578 tems obtain results that are good enough for the application of the system  
579 in a production setting. Specifically, we show that Language Model tuning is  
580 a key factor, as a Language Model-based transfer learning provides the best  
581 performance. The paper has also studied the generalization ability of mod-  
582 els trained in different hospitals, showing that different section types have  
583 significant differences in some cases.

584 The developed automatic annotation models and software are freely avail-  
585 able contacting the authors.

## 586 **Acknowledgements**

587 We gratefully acknowledge the support of NVIDIA Corporation with the  
588 donation of the Titan X Pascal GPU used for this research. This work was  
589 partially funded by the Spanish Ministry of Science and Innovation (DOTT-  
590 HEALTH/PAT-MED PID2019-106942RB-C31), the European Commission  
591 (FEDER), the Basque Government (IXA IT-1343-19), and the EU ERA-  
592 Net CHIST-ERA and the Spanish Research Agency (ANTIDOTE PCI2020-  
593 120717-2).

## 594 **References**

- 595 [1] C. Peterson, C. Hamilton, P. Hasvold, From innovation to implementa-  
596 tion – eHealth in the WHO European Region, World Health Organiza-  
597 tion, 2016.

- 598 [2] M. Adnan, J. Warren, M. Orr, A. Ewens, J. Scott, S. Trubshaw, The  
599 quality of electronic discharge summaries for post-discharge care: Hos-  
600 pital panel assessment and it to support improvement, Health Care and  
601 Informatics Review Online 15 (2011).
- 602 [3] Health at a Glance: Europe 2018 STATE OF HEALTH IN THE EU CY-  
603 CLE, [https://ec.europa.eu/health/sites/default/files/state/  
604 docs/2018\\_healthatglance\\_rep\\_en.pdf](https://ec.europa.eu/health/sites/default/files/state/docs/2018_healthatglance_rep_en.pdf), 2020. Last Online; accessed  
605 31-05-2021.
- 606 [4] Health IT Data Summaries, [https://dashboard.healthit.gov/apps/  
607 health-information-technology-data-summaries.php](https://dashboard.healthit.gov/apps/health-information-technology-data-summaries.php), 2021. Last  
608 Online; accessed 31-05-2021.
- 609 [5] Connecting health and care for the nation: A shared nationwide in-  
610 teroperability roadmap. Office of the National Coordinator for Health  
611 Information Technology (ONC). Washington, DC: U.S. Department of  
612 Health and Human Services (HHS), 2015.
- 613 [6] Recommendation on a European Electronic Health Record exchange  
614 format. European Commision, 2019.
- 615 [7] State of Interoperability among U.S. Non-federal Acute Care Hospitals  
616 in 2018 , [https://www.healthit.gov/sites/default/files/page/2020-  
618 Hospitals-in-2018.pdf](https://www.healthit.gov/sites/default/files/page/2020-03/State-of-Interoperability-among-US-Non-federal-Acute-Care-<br/>617 Hospitals-in-2018.pdf), 2020. Last Online; accessed 31-05-2021.
- 619 [8] openehr, <https://www.openehr.org>, 2020. Last Online; accessed 31-  
620 05-2021.

- 621 [9] Health Level Seven (HL7). FHIR, <http://www.hl7.org>, 2019. Last On-  
622 line; accessed 31-05-2021.
- 623 [10] Health Level Seven (HL7). CDA, <http://www.hl7.org>, 2019. Last On-  
624 line; accessed 31-05-2021.
- 625 [11] H. K., K. Saranto, N. P., Definition, structure, content, use and impacts  
626 of electronic health records: a review of the research literature, *Int J*  
627 *Med Inform.* 77 (5) (2008) 291–304.
- 628 [12] W. LL., Medical records that guide and teach, *N Engl J Med.* 14) (1968)  
629 593–600.
- 630 [13] A. Pomares-Quimbaya, M. Kreuzthaler, S. Schulz, Current approaches  
631 to identify sections within clinical narratives from electronic health  
632 records: a systematic review, *BMC Medical Research Methodology* 19  
633 (2019).
- 634 [14] T. Edinger, D. Demner-Fushman, A. Cohen, S. Bedrick, H. W., Evalu-  
635 ation of Clinical Text Segmentation to Facilitate Cohort Retrieval, in:  
636 *AMIA Annu Symp Proc.*, pp. 660–669.
- 637 [15] J. Lei, B. Tang, X. Lu, K. Gao, M. Jiang, H. Xu, A comprehensive  
638 study of named entity recognition in chinese clinical text, *Journal of the*  
639 *American Medical Informatics Association* 21 (2014).
- 640 [16] Y. Wang, L. Wang, M. Rastegar-Mojarad, S. Moon, F. Shen, N. Afzal,  
641 S. Liu, Y. Zeng, S. Mehrabi, S. Sohn, H. Liu, Clinical information extrac-  
642 tion applications: A literature review, *Journal of Biomedical Informatics*  
643 77 (2018) 34 – 49.

- 644 [17] H.-J. Lee, Y. Zhang, M. Jiang, J. Xu, C. Tao, H. Xu, Identifying direct  
645 temporal relations between time and events from clinical notes, *BMC*  
646 *Medical Informatics and Decision Making* 18 (2018).
- 647 [18] A. Pérez, K. Gojenola, A. Casillas, M. Oronoz, A. Díaz de Ilarraza,  
648 Computer aided classification of diagnostic terms in spanish, *Expert*  
649 *Systems with Applications*, **42**(6), 2949–295 (2015).
- 650 [19] A. Atutxa, A. D. de Ilarraza, K. Gojenola, M. Oronoz, O. P. de Viñaspre,  
651 Interpretable deep learning to map diagnostic texts to icd-10 codes, *In-*  
652 *ternational Journal of Medical Informatics* 129 (2019) 49 – 59.
- 653 [20] K. Xu, M. Lam, J. Pang, X. Gao, C. Band, P. Mathur, F. Papay, A. K.  
654 Khanna, J. B. Cywinski, K. Maheshwari, P. Xie, E. P. Xing, Multimodal  
655 machine learning for automated icd coding, in: F. Doshi-Velez, J. Fack-  
656 ler, K. Jung, D. Kale, R. Ranganath, B. Wallace, J. Wiens (Eds.), *Pro-*  
657 *ceedings of the 4th Machine Learning for Healthcare Conference*, volume  
658 106 of *Proceedings of Machine Learning Research*, PMLR, Ann Arbor,  
659 Michigan, 2019, pp. 197–215.
- 660 [21] A. Duque, H. Fabregat, L. Araujo, J. MartinezRomo, A keyphrasebased  
661 approach for interpretable ICD-10 code classification of Spanish medical  
662 reports, *Artificial Intelligence in Medicine* (2020).
- 663 [22] S. Arnold, R. Schneider, P. Cudré-Mauroux, F. A. Gers, A. Löser, Sec-  
664 tor: A neural model for coherent topic segmentation and classification,  
665 *Transactions of the Association for Computational Linguistics* 7 (2019)  
666 169–184.

- 667 [23] E. Choi, Z. Xu, Y. Li, M. W. Dusenberry, G. Flores, E. Xue, A. M.  
668 Dai, Learning the graphical structure of electronic health records with  
669 graph convolutional transformer, in: Association for the Advancement  
670 of Artificial Intelligence (AAAI).
- 671 [24] S. Rosenthal, K. Barker, Z. Liang, Leveraging medical literature for  
672 section prediction in electronic health records, in: Proceedings of the  
673 2019 Conference on Empirical Methods in Natural Language Processing  
674 and the 9th International Joint Conference on Natural Language Pro-  
675 cessing (EMNLP-IJCNLP), Association for Computational Linguistics,  
676 Hong Kong, China, 2019, pp. 4864–4873.
- 677 [25] E. Rush, I. Danciu, G. Ostrouchov, K. Cho, B. Mayer, Y.-L. Ho, J. Hon-  
678 erlaw, L. Costa, F. Linares, E. Begoli, Jsonize: A scalable machine  
679 learning pipeline to model medical notes as semi-structured documents,  
680 AMIA Joint Summits on Translational Science proceedings. AMIA Joint  
681 Summits on Translational Science 2020 (2020) 533–541.
- 682 [26] L. K. Branting, C. Pfeifer, B. Brown, L. Ferro, J. Aberdeen, B. Weiss,  
683 M. Pfaff, B. Liao, Scalable and explainable legal prediction, Artificial  
684 Intelligence and Law (2020).
- 685 [27] A. R. Terroba, Mejora de la calidad del informe clínico de alta hospita-  
686 laria desde el punto de vista lingüístico, PhD Thesis, University of La  
687 Rioja (2018).
- 688 [28] T. Mikolov, K. Chen, G. S. Corrado, J. Dean, Efficient Estimation of  
689 Word Representations in Vector Space, CoRR abs/1301.3781 (2013).



- 690 [29] J. Pennington, R. Socher, C. D. Manning, Glove: Global vectors for  
691 word representation, in: Empirical Methods in Natural Language Pro-  
692 cessing (EMNLP), pp. 1532–1543.
- 693 [30] T. Mikolov, E. Grave, P. Bojanowski, C. Puhersch, A. Joulin, Advances  
694 in pre-training distributed word representations, in: Proceedings of the  
695 International Conference on Language Resources and Evaluation (LREC  
696 2018).
- 697 [31] I. Yamada, A. Asai, H. Shindo, H. Takeda, Y. Takefuji, Wikipedia2vec:  
698 An optimized tool for learning embeddings of words and entities from  
699 wikipedia, arXiv preprint 1812.06280 (2018).
- 700 [32] A. Akbik, D. Blythe, R. Vollgraf, Contextual string embeddings for  
701 sequence labeling, in: Proceedings of the 27th International Conference  
702 on Computational Linguistics, pp. 1638–1649.
- 703 [33] A. McCallum, W. Li, Early results for named entity recognition with  
704 conditional random fields, feature induction and web-enhanced lexicons,  
705 in: Proceedings of the Seventh Conference on Natural Language Learn-  
706 ing at HLT-NAACL 2003 - Volume 4, CONLL '03, Association for Com-  
707 putational Linguistics, Stroudsburg, PA, USA, 2003, pp. 188–191.
- 708 [34] A. Jagannatha, H. Yu, Bidirectional recurrent neural networks for med-  
709 ical event detection in electronic health records, CoRR abs/1606.07953  
710 (2016).
- 711 [35] A. Atutxa, K. Bengoetxea, A. D. de Ilarraza, M. Iruskieta, Towards a  
712 top-down approach for an automatic discourse analysis for basque: Seg-

- 713       mentation and central unit detection tool, PLoS ONE 14(9): e0221639  
714       (2019).
- 715 [36] H.-J. Dai, S. Syed-Abdul, C.-W. Chen, C.-C. Wu, Recognition and  
716       evaluation of clinical section headings in clinical documents using token-  
717       based formulation with conditional random fields, BioMed Research  
718       International 2015 (2015).
- 719 [37] M. Collins, Discriminative training methods for hidden Markov mod-  
720       els: Theory and experiments with perceptron algorithms, in: Proceed-  
721       ings of the 2002 Conference on Empirical Methods in Natural Language  
722       Processing (EMNLP 2002), Association for Computational Linguistics,  
723       2002, pp. 1–8.
- 724 [38] A. Pérez, R. Weegar, A. Casillas, K. Gojenola, M. Oronoz, H. Dalianis,  
725       Semi-supervised medical entity recognition: A study on spanish and  
726       swedish clinical corpora, Journal of Biomedical Informatics 71 (2017).
- 727 [39] A. Akbik, T. Bergmann, D. Blythe, K. Rasul, S. Schweter, R. Vollgraf,  
728       Flair: An easy-to-use framework for state-of-the-art nlp, in: Proceedings  
729       of the 2019 Conference of the North American Chapter of the Associa-  
730       tion for Computational Linguistics (Demonstrations), pp. 54–59.