

Document downloaded from:

<http://hdl.handle.net/10251/37963>

This paper must be cited as:

Ortigosa, N.; Gimenez, VM. (2014). Raw data extraction from electrocardiograms with Portable Document Format. *Computer Methods and Programs in Biomedicine*. 113(1):284-289. doi:10.1016/j.cmpb.2013.09.014.



The final publication is available at

Copyright Elsevier

1 Raw data extraction from electrocardiograms with  
2 Portable Document Format

3 Nuria Ortigosa<sup>a</sup>, Vicente M. Giménez<sup>a</sup>

4 <sup>a</sup> *I.U. Matemática Pura y aplicada, Universidad Politécnica de Valencia*  
5 *Camino de Vera s/n, 46022 Valencia, Spain*  
6 *nuorar@upvnet.upv.es*

---

7 **Abstract**

During the last two decades there has been a thorough research and development of standards and protocols in order to cope with different electrocardiogram formats from heterogeneous acquisition systems. Despite the efforts of public and private consortiums on creating a standardized electrocardiogram (ECG) storage format, there is still not a single one. Indeed, there is also the necessity of access to raw data of the ECGs previously acquired. Most of these documents have been saved as Adobe PDF files, since for medical staff it is an easy format for later visualization. However, this format presents difficulties when trying to access original raw data for subsequent studies and signal analysis. In this manner, this paper presents an application that obtains plain numerical data from ECG files stored with PDF format. Data can also be exported to one of the most common file formats in existence, to

be easily accessed thereafter.

8 *Keywords:* ECG storage formats, Portable Document Format, Scalable  
9 Vector Graphics

---

## 10 **1. Introduction**

11 The cost of an electrocardiograph is low, specially when compared to  
12 other medical equipment (such as, for example, the radiological one). This  
13 is probably why different ECG recording companies have been working with  
14 such a large number of different data formats. During the last 20 years, the  
15 definition and adoption of a single standardized format for ECG recordings  
16 has been promoted for public and private consortiums [1, 2] due to the fact  
17 that interoperability between them could save around 77.8 billion of dollars  
18 per year [3], just in the United States.

19 Therefore, a European funded project called OpenECG [4] started in 2002  
20 to promote the adoption of the SCP-ECG (“Standard Communications Pro-  
21 tocol for computed assisted ECG”) by implementing visualization tools and  
22 interoperability standards to help manufacturers on their implementations  
23 [5].

24 Similarly, in 2001 the United States Food and Drug Administration (FDA)  
25 requested a standard which eased the storage and exchange of ECGs informa-

26 tion, since it received a large number of annotated ECGs collected in a wide  
27 variety of formats. Thus, the HL7 (a not-for-profit international organiza-  
28 tion for sharing health information) developed the Health Level 7-annotated  
29 ECG (HL7-aECG) standard [6], a new XML-based format.

30 Another format to store medical data is DICOM [7], which initially was  
31 developed for medical image storage and, in 1995, finally became a European  
32 standard also used for cardiac and vascular information.

33 Surprisingly, efforts towards a single standardized ECG format are not  
34 only supported by the Standard Development Organizations. In 2003, Philips  
35 Medical Systems published the XML (extensible markup language) schema  
36 that they used for their entire line of ECG products [8], and began to de-  
37 liver this information to its costumers and to European OpenECG project  
38 members [9, 10].

39 Nevertheless, although interoperability farther than at regional or na-  
40 tional level is a desirable goal, it may still take another 20 years to be fully  
41 achieved [11]. In this sense, integrated applications for standardized data  
42 formats exchange have been developed to cope with the numerous different  
43 formats and to facilitate their visualization [12]. For example, van Ettinger  
44 et al. [13] have implemented a conversion library and an ECG viewer to

45 work with HL7-aECG, SCP-ECG and DICOM files. In order to facilitate  
46 interoperability for records obtained in Italy, Marcheschi et al. [14] designed  
47 a network infrastructure to manage with different standards. Similarly, [15]  
48 and [16] use XML language as a central platform to exchange data between  
49 these formats. Moreover, in [17] Trigo et al. have recently developed a  
50 modular application to exchange ECGs information of different data formats  
51 across healthcare information systems.

52 Unfortunately, ECGs previously stored often present data formats dif-  
53 ferent from those commented above. In particular, many electrophysiolo-  
54 gists and medical staff usually save only the plots of the ECG recordings as  
55 Adobe PDF (“Portable Document Format”) due to the fact that, when ECG  
56 acquisition equipment is replaced in hospitals, ECG recordings stored with  
57 proprietary formats hamper the free access to previous ECGs in the med-  
58 ical history. Therefore, PDF format assures them an easy visualization in  
59 any computer and allows subsequent revision for clinical studies and patients  
60 evolution analysis. Thus, in this paper we present an application developed  
61 to recover raw data from ECG recordings saved in PDF format, in order to  
62 facilitate later signal processing of ECGs previously stored by medical staff  
63 in this format, and to integrate them with the existing information system

64 of the hospital.

65 The rest of the paper is organized as follows. Features of ECGs stored  
66 with PDF format, a description of their structure, and the application inter-  
67 face are described in Section 2. Section 3 reports the application computa-  
68 tional requirements, some particularities and discussion of results. Finally,  
69 conclusions and future work are drawn in Section 4.

## 70 **2. Materials and Methods**

### 71 *2.1. SVG files*

72 The ECGs recorded with PDF format whose raw data we wanted to  
73 extract were acquired using the Philips PageWriter TC50 [18]. When revising  
74 the technical sheets of this electrocardiograph we noticed that, when Philips  
75 Medical Systems began to develop their ECG data format based on XML,  
76 they turned to the Scalable Vector Graphics (SVG) application language [9],  
77 which is ideal to display easily two-dimensional graphics.

78 SVG [19] is a markup language which is able to formalize a set of graphical  
79 elements such as rectangles, lines, polygons, etc. The most relevant feature  
80 of SVG documents is that they can be considered simultaneously as both  
81 text and images. This ambivalence points out to SVG as an ideal format to

82 connect textual and graphical information in just one file [20].

83 Therefore, we first tried to convert electrocardiograms to SVG data for-  
84 mat, to check whether the Philips electrocardiograph stored PDF files as  
85 vector graphics. In order to do so, we used the free software tool Inkscape  
86 for conversion [21]. We confirmed that the converted SVG file was formed  
87 by a set of graphical elements located with their absolute coordinates, and  
88 not by a bitmap data set. As a result, we decided to extract raw data by  
89 implementing an application whose input were a PDF file and the generated  
90 output contained the corresponding true-value numerical data of each lead,  
91 in order to make easier the subsequent signal analysis.

## 92 *2.2. Electrocardiogram leads*

93 An electrocardiogram is the main instrument for the diagnosis of cardio-  
94 vascular diseases. It can be defined as the graphical representation of the  
95 electrical activity of the heart.

96 In electrocardiography, the term “lead” refers to the measurement of volt-  
97 age between two electrodes, which are placed on the patient’s body. To per-  
98 form a standard 12-lead ECG, it takes 10 electrodes: 6 electrodes for the  
99 chest leads (V1, V2, V3, V4, V5, V6), and 4 electrodes to acquire the limb  
100 leads (I, II, III) and the augmented limb leads (aVL, aVF, aVR).

101 Depending on the duration of the leads that the electrophysiologist prefers  
102 to visualize, ECGs can be stored using different printout formats. Ones of  
103 the most popular printout formats are:  $3 \times 4$ ,  $3 \times 4 + 1$ , and  $6 \times 2$ . Details  
104 of these lead organization will be provided in Figure 2 and Section 3.

### 105 *2.3. Electrocardiogram structure*

106 As aforementioned, the SVG file obtained by Inkscape from ECG PDF  
107 conversion presents a structure similar to an XML file. The horizontal and  
108 vertical coordinates (in pixels) of each sample from the ECG lead is repre-  
109 sented under a label called  $\langle path \rangle$ . Coordinates can be defined as abso-  
110 lute or relative to the last coordinate, depending on whether after the label  
111  $\langle path \rangle$  the token is  $d = "M$  or  $d = "m$ , respectively. In our case, we have  
112 configured Inkscape options to generate SVG files with absolute coordinates,  
113 being the abscissa and the ordinate coordinates separated by a comma.

114 We must remark that, under the label  $\langle path \rangle$  there appear not only the  
115 polylines which correspond to the ECG leads, but also the lines that define  
116 the background grid and the reference pulses. They can be differentiated  
117 by their corresponding line widths and colours, as well as by their order of  
118 appearance along the SVG file and the number of samples for which they are  
119 defined.



#### 120 2.4. Application interface

121 The presented application has been programmed using GUIDE, the MAT-  
122 LAB graphical user interface development environment. Figure 1 shows the  
123 main window of the ECG data extraction application.

124 The interface is divided in three subsections. In the left top half the input  
125 ECG parameters are shown, such as the paper speed and the amplitude scale  
126 of the limb and the chest leads. It is standard to represent each microVolt  
127 ( $\mu V$ ) of amplitude as 10 mm, and each second as 25 mm. However, as  
128 faster paper speeds and different scales can be used, these parameters can be  
129 modified by the user.

130 Then, the user must proceed to load the desired PDF file whose raw data  
131 is going to be extracted. Once it has been selected, the application calls  
132 the Inkscape program as a background task, in order to perform the PDF  
133 conversion to SVG format. All our ECGs acquired and recorded information  
134 of 12 leads, according to one of the following printout formats:  $3 \times 4$ ,  $3 \times$   
135  $4 + 1$  or  $6 \times 2$  (see Figure 2). So, next step consists of automatically detect  
136 which lead configuration was used when the loaded ECG was recorded. This  
137 automatic recognition takes place by detecting and counting the number of  
138 reference pulses.

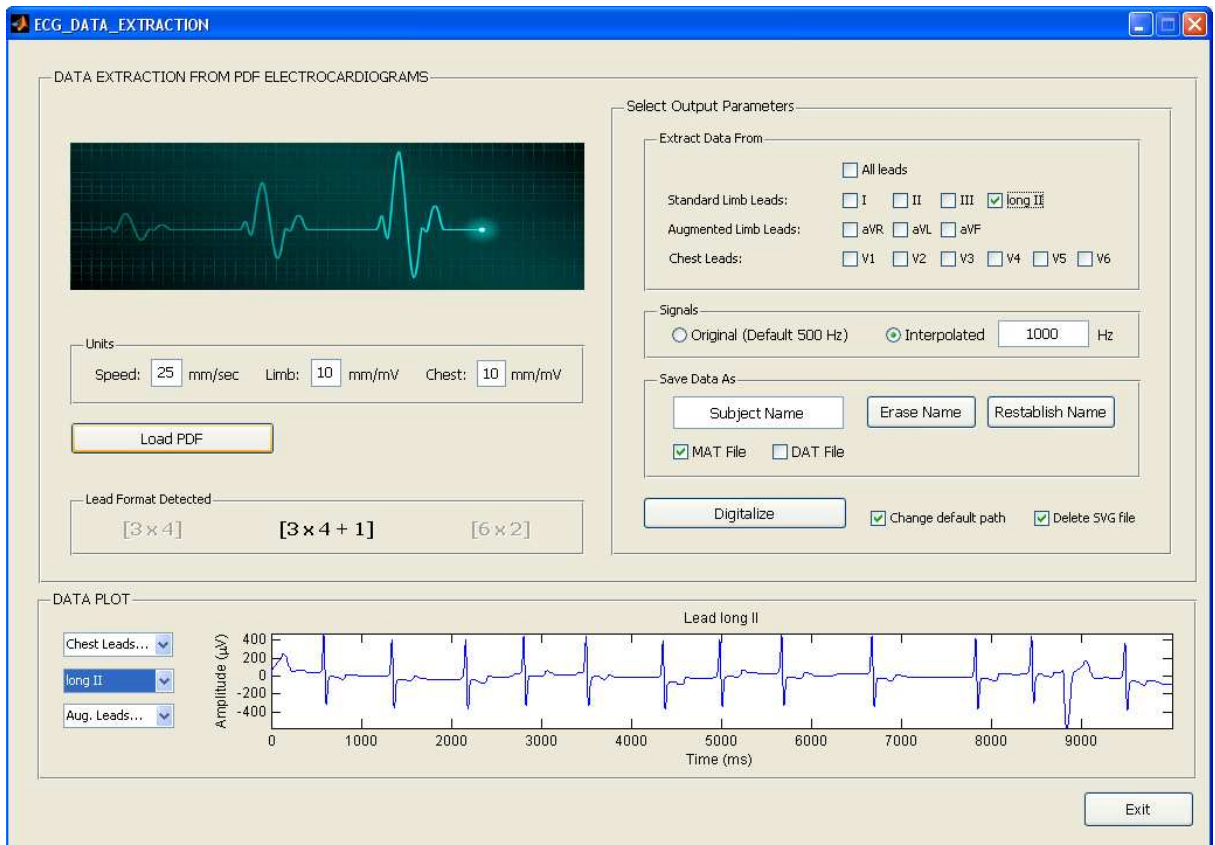
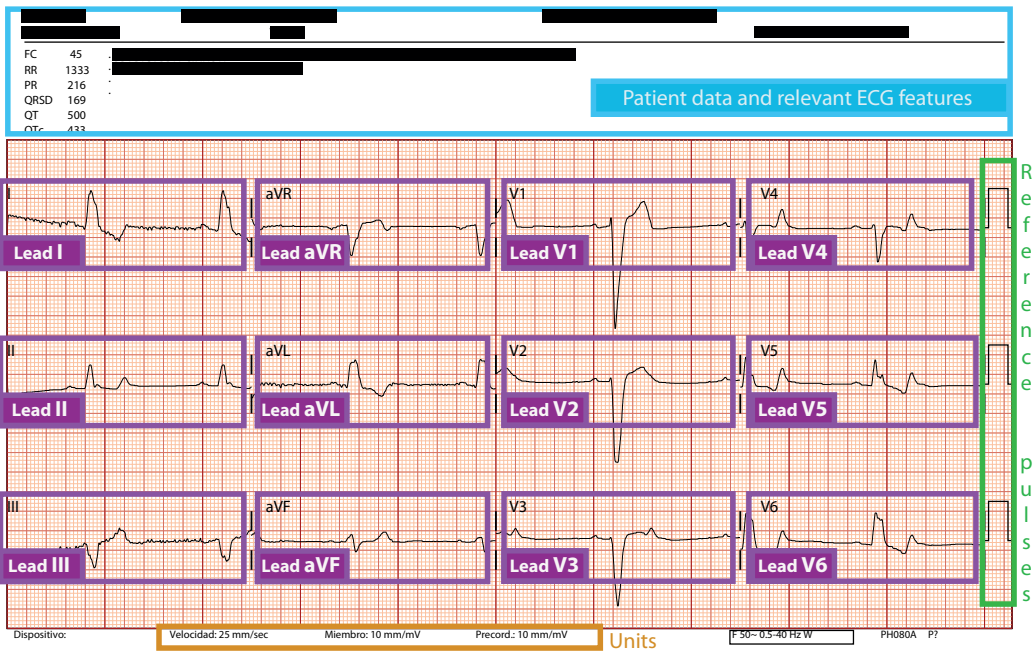
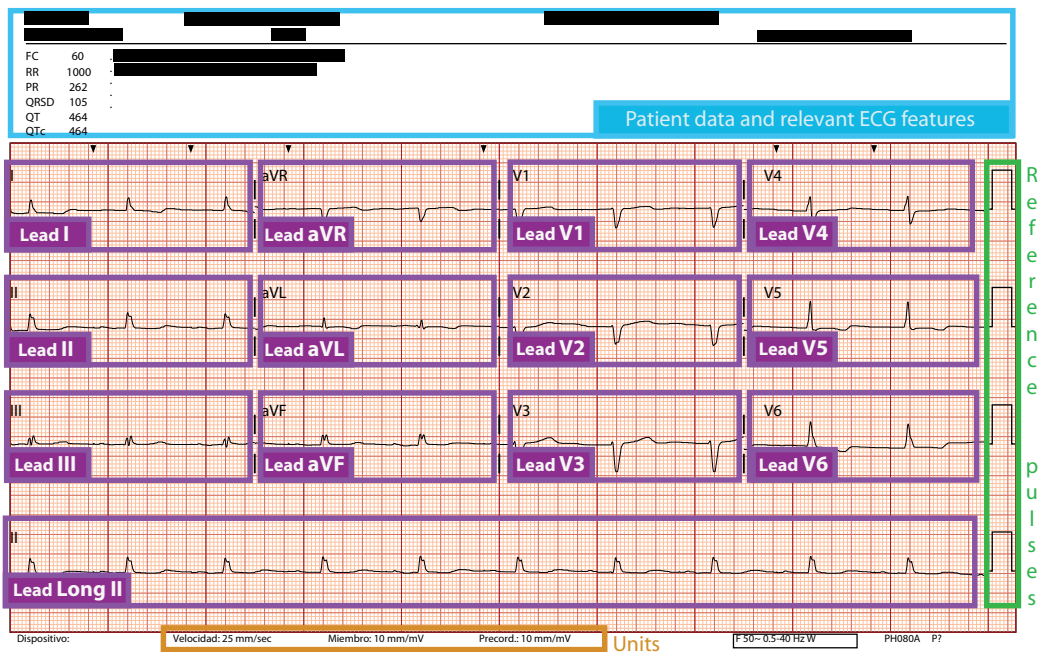


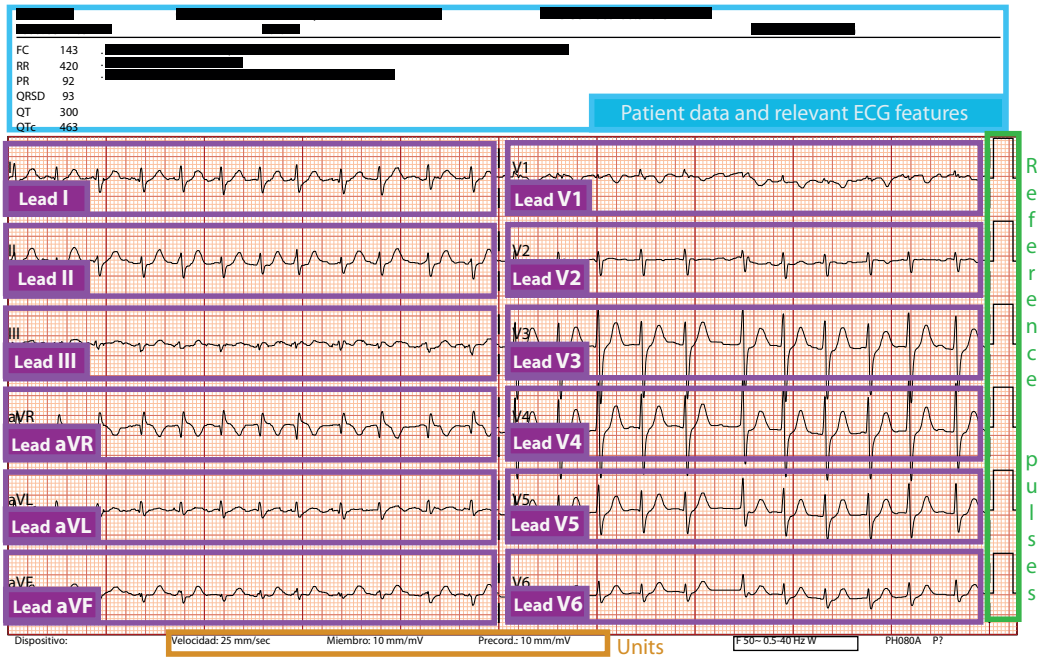
Figure 1: Graphical interface of the presented application for raw data extraction from ECGs in PDF format.



(a)



(b)



(c)

Figure 2: Most common ECG printout formats. Subfigure 2(a) shows  $3 \times 4$  leads configuration, whereas 2(b) and 2(c) depict  $3 \times 4 + 1$  and  $6 \times 2$ , respectively.

139 Next, the SVG file generated by Inkscape is read to look for the numerical  
 140 information of leads, depicted under the label  $\langle path \rangle$  and the token  $d =$   
 141 “ $M$ ”. As commented in section 2.3, every graphical element of the SVG file is  
 142 indicated using those tokens. Consequently, so as to ensure the differentiation  
 143 of the grid from the reference pulses and the leads, the application takes into  
 144 account not only the colour and width of the lines, but also the number of

145 samples which are contained under each label (8 in the case of the reference  
146 pulse, and 2 for each line of the grid). Afterwards, once the numerical data  
147 of each lead has been extracted, the application translates this coordinates  
148 to the origin by subtracting the abscissa and ordinate coordinates of the  
149 corresponding reference pulse for each lead. For example, in case of  $3 \times 4 + 1$   
150 lead configuration, content of leads I, aVR, V1 and V4 are relocated using  
151 the first reference pulse information (Figure 2(b)).

152 Thereafter, based on the information content of the paper speed and the  
153 leads scale textfields, the application converts the extracted coordinates of  
154 each sample (which are expressed in pixels) to true values in milliseconds and  
155  $\mu V$ .

156 At the right top half of the interface, the user can indicate the output  
157 parameters of the extracted data file. In case the user needed the leads raw  
158 data with a certain frequency sampling, he/she could indicate the desired  
159 value in the textfield next to the button *Interpolated*. The user can also  
160 choose the leads he wants to extract, the name of the output file (which is  
161 the same as the PDF file by default), and the directory where it will be saved,  
162 as long as the format to save this data: with .MAT extension (the default  
163 binary format for data files in MATLAB) or with .DAT extension (an ASCII

164 data file whose content can be read just using any text editor application,  
165 such as Wordpad, for example, or be plotted with gnuplot).

166 Finally, at the bottom of the interface, the user can choose and visualize  
167 the data extracted from the desired leads. Figure 3 depicts a flowchart with  
168 the algorithm and the tasks developed by the application.

### 169 **3. Discussion**

170 The application presented in this paper has been developed with the  
171 purpose of extracting numerical information from ECGs that have been pre-  
172 viously recorded as PDF-format files. Although most of our ECGs were  
173 acquired using the Philips PageWriter TC50, the presented application is  
174 able to obtain numerical data from ECGs stored with PDF format using any  
175 model of Philips electrocardiograph, and those from other brands which are  
176 based on XML. This is of great value, since this implies a large proportion of  
177 ECG machines, due to the fact that Philips (via its acquisition of Hewlett-  
178 Packard Medical Products Group) represents one of the largest suppliers of  
179 ECG machines in the world.

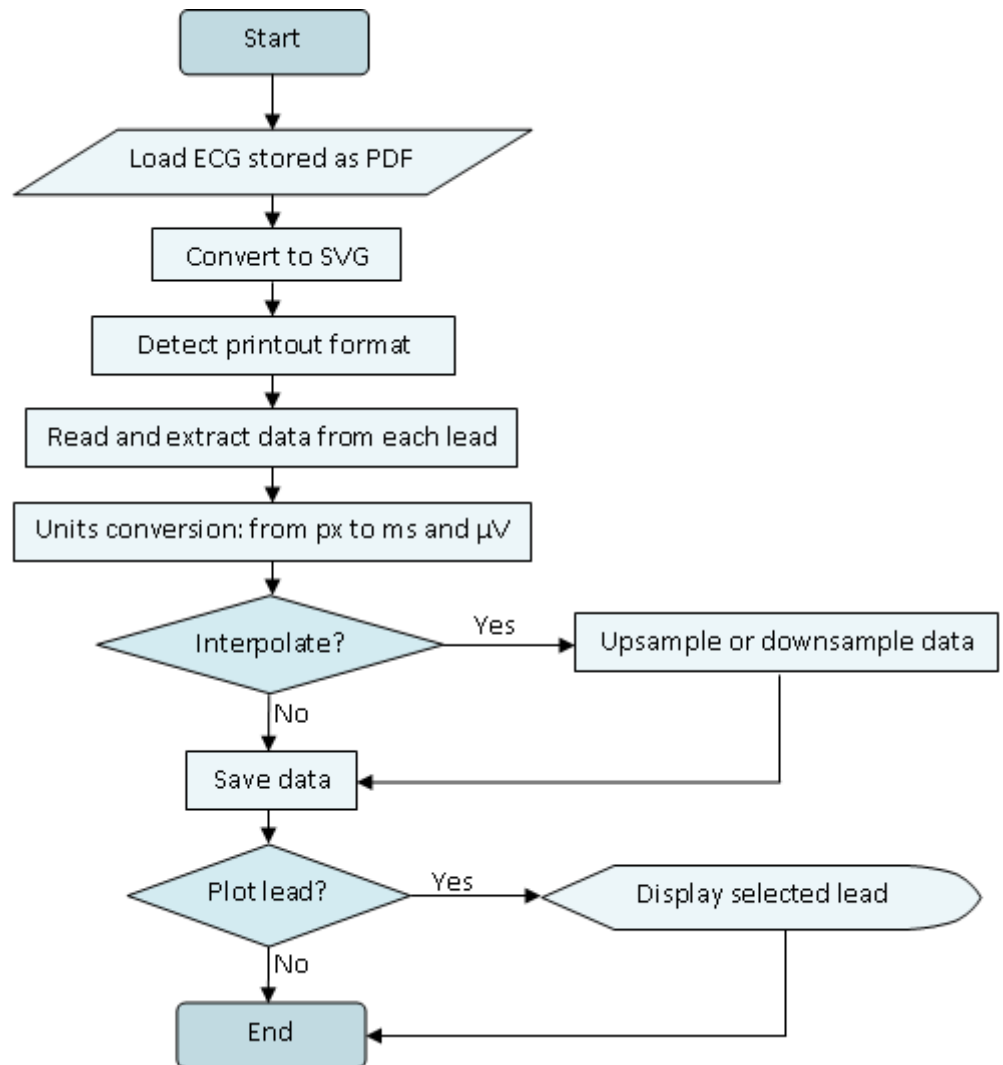


Figure 3: Flowchart of the tasks developed by the application.

180 The valid configurations for automatic recognition of lead organization  
181 are the most common ones:  $3 \times 4$ ,  $3 \times 4 + 1$  or  $6 \times 2$ , as commented in section  
182 2.4. For our ECGs, the first configuration corresponds to all 12 leads of 2.5  
183 seconds of duration. The second one is the same, but adding information  
184 for 10 seconds of lead II (long II). In the last configuration ( $6 \times 2$ ), all 12  
185 leads have a duration of 5 seconds, displayed along 6 rows at the PDF. In  
186 case we wanted to include any other different configuration of leads, software  
187 modifications we should do consist of looking for reference pulses and relocate  
188 the origin of associated leads according to them.

189 Regarding the detection of relevant marks in the ECG (such as those  
190 produced by pacemakers or indicated by the cardiologist, as shown in Figure  
191 4), it is important to remark that they also appear in the SVG file under  
192 the label  $\langle path \rangle$ . Thus, in case they could be of interest for subsequent  
193 analysis, we could extract their temporal position, just looking for polylines  
194 defined by three coordinates with filled style (which define the triangle of the  
195 mark).



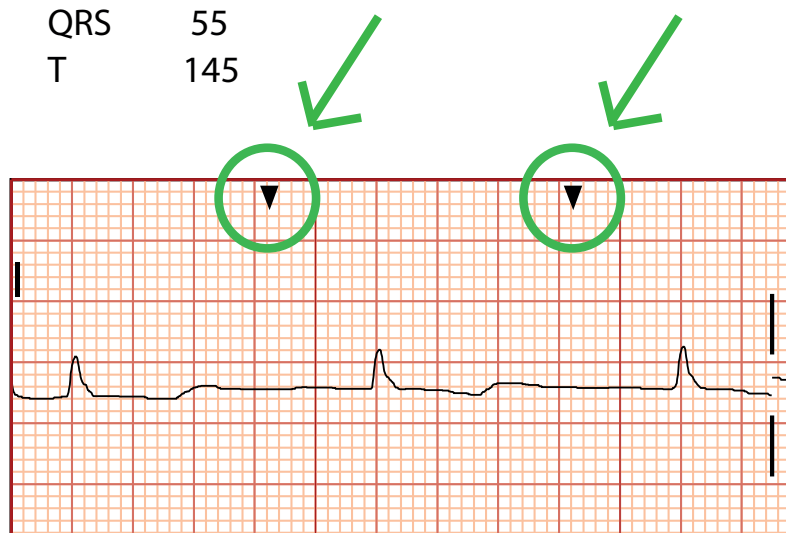


Figure 4: Example of ECG with additional marks.

196        Concerning the sampling frequency of the raw data, we should take into  
197 account that, from SVG files converted from PDF format, we can extract  
198 500 samples per second. However, these data are not uniformly sampled due  
199 to the process of storage of the data. In order to to facilitate the subsequent  
200 analysis we provide an output datafile with uniformly interpolated data. In  
201 addition, in case we needed a different sampling frequency, we could indicate  
202 the desired one at the corresponding text box, and the generated output file  
203 will be saved by uniformly upsampling or downsampling the original data.  
204 Regarding uniformly sampled data reconstruction, we tested with different  
205 interpolation methods, such as nearest-neighbour, linear, cubic and splines.

206 We finally chose linear interpolation method since it presents good results  
207 and low computational cost.

208 Besides, it is important to remind that the presented application carries  
209 out the conversion to SVG format by means of the open source software  
210 Inkscape. In order to ensure the proper working, the application checks all  
211 the existing hard disk drives, so as to find the “Program Files\Inkscape” di-  
212 rectory. Moreover, it also differentiates when using a 64 bits operating system  
213 or a 32 bits one, by looking for the directory “Program Files (x86)\Inkscape”  
214 in that case.

215 Finally, we are going to analyse the computational load of the whole data  
216 recovering for our application. As aforesaid, it has been implemented by  
217 using the MATLAB graphical user interface development environment. De-  
218 tailed computation times spent under a 2.79GHz CPU with 3GB of RAM  
219 are shown in Table 1. It can be seen that for lead configuration  $6 \times 2$ , the  
220 program spends larger times when extracting data and saving the raw data  
221 output file. It is due to the larger number of samples of this configuration  
222 (30000 when sampling frequency is 500Hz) compared to other printout for-  
223 mats (15000 and 20000 samples for  $3 \times 4$  and  $3 \times 4 + 1$ , respectively). Anyway,  
224 on average, computation times are lower than 10 seconds, and they could be

225 much more smaller (less than 4 seconds per ECG) in case of minimal inter-  
 226 face interaction; indeed, the program can be run in batch mode for which  
 227 the user must indicate just the directory where the ECG PDF-files are stored  
 228 and the output directory where the transformed files will be saved.

Table 1: Computation times (in seconds) for different tasks of the presented application for a CPU with 3GB of RAM, a 2.79GHz processor and a Windows XP operating system. First column corresponds to time spent on calls to Inkscape (which represents a quarter of the time).

Printout format	Load PDF and SVG conversion	Lead configuration recognition	Data extraction and output data file creation
$3 \times 4$	2.216	2.223	3.064
$3 \times 4 + 1$	1.737	2.221	5.860
$6 \times 2$	2.080	2.427	8.586

## 229 4. Conclusions

230 In this paper we have presented an application that extracts raw data  
 231 from ECG stored with PDF format. It is based on the conversion of PDF  
 232 files to an XML based data format, the Scalable Vector Graphics (SVG),  
 233 which describes two-dimensional graphical objects. Thus, the application is

234 able to obtain pixel coordinates of the polylines that conform the different  
235 leads of the ECG, so as to convert them to amplitudes of the ECG signal  
236 (in  $\mu V$ ). The application can also automatically recognize the most common  
237 ECG leads' configuration and can provide an ASCII output data file with  
238 different sampling frequencies.

239 As aforementioned, the application is very useful to recover data from  
240 stored ECGs with PDF format. Cardiologists and medical staff usually save  
241 the ECG recordings with this format, so as to ensure the later visualization  
242 without problems of format compatibility when using proprietary data for-  
243 mats from different ECG devices. The presented application recovers data  
244 from these files, allowing to perform numerical analysis and subsequent stud-  
245 ies of patient medical evolution along time.

246 Future work will include the conversion of recovered data from PDF ECGs  
247 to standard formats (such as SCP, for example), which would ease the patient  
248 history integration of recent ECG recordings with those acquired previously,  
249 favouring the adoption of standards and interoperability.

250

## 251 **Appendix**

252       The application has been developed using MATLAB software, and its ex-  
253       ecutable standalone application for Microsoft Windows Operating Systems is  
254       available at <http://personales.upv.es/nuorar/>. It runs on any Microsoft Win-  
255       dows Operating System (from XP and later), when the MATLAB Compiler  
256       Runtime (MCR) library is installed on the computer.

257       The user should download and install the open source Inkscape vector  
258       graphics editor and the MCR library before running the .exe file. Instructions  
259       for installation are detailed at README.pdf file at the above-mentioned  
260       website.

261       We also plan to export this software to Unix and OSX operating systems.  
262       As soon as available, these standalone applications will be uploaded to the  
263       above-mentioned website to be fully accessible.

## 264 **Acknowledgements**

265       This work was supported by grants MTM2010-15200 (from the Spanish  
266       Ministry of Economy and Competitiveness) and UPV-IIS La Fe, 2012/0468.

267 **References**

- 268 [1] R.R. Bond, D.D. Finlay, C.D. Nugent, and G. Moore. A review of  
269 ECG storage formats. *International Journal of Medical Informatics*,  
270 80(10):681–697, 2011.
- 271 [2] J.D. Trigo, A. Alesanco, I. Martínez, and J. García. A review on digital  
272 ECG formats and the relationships between them. *IEEE Transactions*  
273 *on Information Technology in Biomedicine*, 16(3):432–444, May 2012.
- 274 [3] J. Walker, E. Pan, D. Johnston, J. Adler-Milstein, D.W. Bates, and  
275 B. Middleton. The value of health care information exchange and inter-  
276 operability. *Health Affairs*, pages W5–10–W5–18, January 2005.
- 277 [4] OpenECG project. <http://www.openecg.net/>, Last accessed on June  
278 20, 2013.
- 279 [5] C.E. Chronaki, F. Chiarugi, P.J. Lees, M. Bruun-Rasmussen, F. Con-  
280 forti, R. Ruiz-Fernández, and C. Zywieta. OpenECG: An european  
281 project to promote the SCP-ECG standard, a further step towards in-  
282 teroperability in electrocardiography. pages 285–288. *Computers in Car-*  
283 *diology*, 2002.

- 284 [6] B.D. Brown and F. Badilini. HL7 aECG implementation guide.  
285 [http://www.hl7.org/implement/standards/product\\_brief.cfm?product\\_id](http://www.hl7.org/implement/standards/product_brief.cfm?product_id=102)  
286 [=102](http://www.hl7.org/implement/standards/product_brief.cfm?product_id=102), Last accessed on June 20, 2013.
- 287 [7] DICOM (digital imaging and communications in medicine).  
288 <http://medical.nema.org/>, Last accessed on June 20, 2013.
- 289 [8] N. Long. Open ECG data standard: Philips medical systems perspec-  
290 tive. *Journal of Electrocardiology*, 36:167, 2003.
- 291 [9] E.D. Helfenbein, R. Gregg, and S. Zhou. Philips medical systems sup-  
292 port for open ECG and standardization efforts. pages 393–396. *Com-*  
293 *puters in Cardiology*, 2004.
- 294 [10] S. Zhou and E. Helfenbein. OpenECG format: Philips’ experience. pages  
295 46–47. 2nd OpenECG Workshop, 2004.
- 296 [11] Semantic interoperability for better health and safer healthcare.  
297 [http://www.empirica.com/publikationen/documents/2009/semantic-](http://www.empirica.com/publikationen/documents/2009/semantic-health-report.pdf)  
298 [health-report.pdf](http://www.empirica.com/publikationen/documents/2009/semantic-health-report.pdf), Last accessed on June 20, 2013.
- 299 [12] C.E. Chronaki, F. Chiarugi, A. Macerata, F. Conforti, H. Voss, I. Jo-  
300 hansen, R. Ruiz-Fernández, and C. Zywietz. Interoperability in digital

- 301 electrocardiography after the openECG project. pages 49–52. Comput-  
302 ers in Cardiology, 2004.
- 303 [13] M.J.B. van Ettinger, J.A. Lipton, M.C.J. de Wijs, N. van der Putten,  
304 and S.P. Nelwan. An open source ECG toolkit with DICOM. pages  
305 441–444. Computers in Cardiology, 2008.
- 306 [14] P. Marcheschi, A. Mazzarisi, S. Dalmiani, and A. Benassi. ECG stan-  
307 dards for the interoperability in patient electronic health records in italy.  
308 pages 549–52. Computers in Cardiology, 2006.
- 309 [15] X. Li, V. Vojisavljevic, and Q. Fang. An XML based middleware for  
310 ECG format conversion. pages 1691–1694. 31st Annual International  
311 Conference of the IEEE EMBS, 2009.
- 312 [16] J.D. Trigo, A. Kollmann, A. González, D. Hayn, A. Alesanco,  
313 G. Schreier, and J. García. Plataforma para la integración y la gestión  
314 homogénea de formatos de electrocardiografía. XXVIII Congreso Anual  
315 de la Sociedad Española de Ingeniería Biomédica CASEIB, 2010.
- 316 [17] J.D. Trigo, I. Martínez, A. Alesanco, A. Kollmann, J. Escayola, D. Hayn,  
317 G. Schreier, and J. García. An integrated healthcare information system  
318 for end-to-end standardized exchange and homogeneous management of



- 319 digital ECG formats. *IEEE Transactions on Information Technology in*  
320 *Biomedicine*, 16(4):518–529, July 2012.
- 321 [18] Philips pagewriter TC50 cardiograph. [http://www.healthcare.](http://www.healthcare.philips.com/main/products/cardiography/products/cardiograph/pagewritertc50.wpd)  
322 [philips.com/main/products/cardiography/products/cardiograph/pagewritertc50.wpd](http://www.healthcare.philips.com/main/products/cardiography/products/cardiograph/pagewritertc50.wpd), Last accessed on June 20, 2013.
- 324 [19] Scalable vector graphics (SVG). <http://www.w3.org/TR/SVG/>, Last  
325 accessed on June 20, 2013.
- 326 [20] A. de la Rosa and J.A. Senso. La dualidad texto-imagen en SVG (scal-  
327 able vector graphics): nuevas posibilidades para la descripción de in-  
328 formación gráfica. *El Profesional de la Información*, 12(5):377–398,  
329 September 2003.
- 330 [21] Inkscape: Open source scalable vector graphics editor. [http://inkscape](http://inkscape.org/)  
331 [.org/](http://inkscape.org/), Last accessed on June 20, 2013.