

bioPDFX: preparing PDF scientific articles for biomedical text mining

Shitij Bhargava¹, Tsung-Ting Kuo², Ankit Goyal¹, Vincent Kuri¹, Gordon Lin¹, Chun-Nan Hsu^{Corresp.}²

¹ Department of Computer Science and Engineering, Jacobs School of Engineering, University of California, San Diego, La Jolla, California, United States

² Health System Department of Biomedical Informatics, School of Medicine, University of California, San Diego, La Jolla, California, United States

Corresponding Author: Chun-Nan Hsu
Email address: chunnan@ucsd.edu

Background. There is huge amount of full-text biomedical literatures available in public repositories like PubMed Central (PMC). However, a substantial number of the papers are in Portable Document Format (PDF) and do not provide plain text format ready for text mining and natural language processing (NLP). Although there exist many PDF-to-text converters, they still suffer from several challenges while processing biomedical PDFs, such as the correct transcription of titles/abstracts, segmenting references/acknowledgements, special characters, jumbling errors (the wrong order of the text), and word boundaries.

Methods. In this paper, we present bioPDFX, a novel tool which complements weaknesses with strengths of multiple state-of-the-art methods and then applies machine learning methods to address all issues above

Results. The experiment results on publications of Genome Wide Association Studies (GWAS) demonstrated that bioPDFX significantly improved the quality of XML comparing to state-of-the-art PDF-to-XML converter, leading to a biomedical database more suitable for text mining.

Discussion. Overall, the whole pipeline developed in this paper makes the published literature in form of PDF files much better suited for text mining tasks, while slightly improving the overall text quality as well. The service is open to access freely at URL: <http://textmining.ucsd.edu:9000>. A list of PubMed Central IDs of the 941 articles (see Supplemental File 1) used in this study is available for download at the same URL. The instructions of how to run the service with a PubMed ID are described in Supplemental File 2.

1 **bioPDFX: preparing PDF scientific articles for biomedical**
2 **text mining**

3
4 Shitij Bhargava¹, Tsung-Ting Kuo², Ankit Goyal¹, Vincent Kuri¹, Gordon Lin¹, Chun-Nan Hsu²

5
6 ¹Department of Computer Science and Engineering, Jacobs School of Engineering,
7 University of California, San Diego, La Jolla, California, United States

8 ²Health System Department of Biomedical Informatics, School of Medicine,
9 University of California, San Diego, La Jolla, California, United States

10

11 Corresponding Author:

12 Chun-Nan Hsu

13 9500 Gilman Drive, La Jolla, CA 92093, United States

14 Email address: chunnan@ucsd.edu

15

16 **Abstract**

17 **Background.** There is huge amount of full-text biomedical literatures available in public
18 repositories like PubMed Central (PMC). However, a substantial number of the papers are in
19 Portable Document Format (PDF) and do not provide plain text format ready for text mining and
20 natural language processing (NLP). Although there exist many PDF-to-text converters, they still
21 suffer from several challenges while processing biomedical PDFs, such as the correct transcription
22 of titles/abstracts, segmenting references/acknowledgements, special characters, jumbling errors
23 (the wrong order of the text), and word boundaries.

24
25 **Methods.** In this paper, we present bioPDFX, a novel tool which complements weaknesses with
26 strengths of multiple state-of-the-art methods and then applies machine learning methods to
27 address all issues above.

28
29 **Results.** The experiment results on publications of Genome Wide Association Studies (GWAS)
30 demonstrated that bioPDFX significantly improved the quality of XML comparing to state-of-the-
31 art PDF-to-XML converter, leading to a biomedical database more suitable for text mining.

32
33 **Discussion.** Overall, the whole pipeline developed in this paper makes the published literature in
34 form of PDF files much better suited for text mining tasks, while slightly improving the overall
35 text quality as well. The service is open to access freely at URL: <http://textmining.ucsd.edu:9000>.
36 A list of PubMed Central IDs of the 941 articles (see Supplemental File 1) used in this study is
37 available for download at the same URL. The instructions of how to run the service with a PubMed
38 ID are described in Supplemental File 2.

40 1. Introduction

41 PubMed Central (PMC) (2015e) contains full-text of about 4 million biomedical literatures,
42 and is one of the most important freely available data sources for the biomedical natural language
43 processing (NLP) research field. Although many new publications now provide plain text along
44 with the Portable Document Format (PDF), a substantial number of them do not. However, text
45 mining algorithms work more effectively on text-based formats such as plain text, XML or HTML
46 documents. Therefore, NLP researchers usually need to transcribe the biomedical literatures from
47 PDF to text as a preprocessing step of the NLP pipeline.

48 However, transcribing PDF to text accurately is not trivial. The design goals of the PDF
49 standard target ease of human readability rather than electronic consumption of data by other
50 software tools. That is, the PDF standard does not attempt to encode any semantic connection
51 between characters in a word or between paragraphs, but characters are "painted" individually at
52 specific 2-dimensional coordinates. Berg et. al. outlined some challenges encountered in
53 transcribing PDF documents to text accurately, such as varied reading orders, sectional formatting
54 and number of columns (Berg 2011). These issues create a variability in text quality in the
55 generated outputs and impedes accurate text mining.

56 Although there are already a variety of tools available today for PDF-to-text transcription
57 (such as Apache PDFBox Mozilla PDF.js (2015d) Adobe Acrobat SDK (2015a) Tesseract
58 Optical Character Recognition (OCR) (Smith 2007) (Constantin et al. 2013) several
59 challenges to be solved for converting biomedical literature PDF to XML for text mining purpose::

- 60 1. *Extracting title and abstracts correctly.* In case of manuscripts which might have contents
61 other than title or abstract on their first page, conversion tools like PDFX struggle to

62 identify them correctly.

63 2. *Identifying reference and acknowledgement sections in various format.* The format of
64 references varies considerably among different publications ranging from a separate
65 section towards the end to individual references in footnotes continuing alongside the
66 content body. Take PDFX as an example, as a purely rule-based system, it often falters in
67 recognizing references when the common pattern of having references in the end with a
68 heading is not present. Similar problems were noticed in identifying acknowledgements as
69 well.

70 3. *Recognizing special characters (such as '!' and '@').* PDFX and all PDF to text converters
71 usually suffer from the problem of recognizing the special characters depending on the
72 publishing of the PDF files.

73 4. *Detecting and fixing jumbling errors (the wrong order of text while converting PDF to*
74 *text).* Often most PDF conversion tools jumbles the work order in or across sentences due
75 to mistakes in reading order detection/segmentation. An example is given in Figure 1 where
76 the exponents of numbers in a table are extracted as a separate column and hence are jumbled
77 with respect to the actual text.

78 5. *Correcting word boundaries.* Often words straddling the column boundary have a hyphen
79 inserted in between them (e.g., “cancer” might become “can-” and “cer” if it is the last
80 word of a row). Although this is not an error due to conversion, we would still like to
81 correct it and merge the hyphenated parts as a single word wherever possible to support
82 NLP afterwards.

83

84 **Figure 1. A kind of jumbling error in a table.**

85 The table at the top is what is present in the PDF. The lower portion shows the transcribed XML.

86 The exponents marked in the red-colored box are incorrectly extracted as a separate column in the

87 XML.

A. Discovery group.

dbSNP	Chromosome	Physical Position	Gene Relationship	Gene Distance	Gene	BP Change	Prob Allele	Prob Genotype
rs6975107	7	120168143	intron	0	KCND2	C-> T	4.15x10 ⁻⁰⁹	3.74x10 ⁻⁰⁹
rs318125	X	97090551	downstream	348,298	DIAPH2	C-> T	4.35x10 ⁻⁰⁹	1.25x10 ⁻⁰⁵
rs5916727	X	104294270	intron	0	IL1RAPL2	C-> T	6.66x10 ⁻⁰⁹	1.34x10 ⁻⁰⁵
rs11863929	16	86861934	upstream	159,445	ZNF469	C-> G	1.77x10 ⁻⁰⁸	5.48x10 ⁻⁰⁷
rs1578826	X	85889112	intron	0	DACH2	T-> C	2.20x10 ⁻⁰⁸	5.46x10 ⁻⁰⁶
rs7616661	3	5940543	downstream	703,894	EDEM1	T-> G	4.82x10 ⁻⁰⁸	3.37x10 ⁻⁰⁸

<region class="DoCO:TextChunk" id="586" confidence="possible" page="26"

column="1">A. Discovery group.</region>

<outsider class="DoCO:TextBlock" type="sidenote" id="587" page="26"

column="1">-09 -05 -05 -07 -06 -08</outsider>

-<region class="unknown" id="588" page="26" column="1">

-09 rs6975107 7 120168143 intron 0 KCND2 C-> T 4.15x10 3.74x10 -09 rs318125 X
97090551 downstream 348,298 DIAPH2 C-> T 4.35x10 1.25x10 -09 rs5916727 X
104294270 intron 0 IL1RAPL2 C-> T 6.66x10 1.34x10 -08 rs11863929 16 86861934
upstream 159,445 ZNF469 C-> G 1.77x10 5.48x10 -08 rs1578826 X 85889112 intron 0
DACH2 T-> C 2.20x10 5.46x10 -08 rs7616661 3 5940543 downstream 703,894 EDEM1
T-> G 4.82x10 3.37x10

88 </region>

89

90 To address these issues, we proposed bioPDFX, a novel tool that integrates current state-of-
91 the-art PDF conversion methods to transcribe biomedical PDF articles to high quality text in XML
92 format. The bioPDFX tool leverages the following four tools to transcribe text of the given PDF:
93 *PDFX* (Constantin et al. 2013), *PMC Entrez e-utilities API* (Sayers et al. 2011), *Tesseract OCR*
94 (Smith 2007), and *Apache PDFBox* (2015b). The input of bioPDFX includes a biomedical
95 literature PDF file and its corresponding PubMed ID, while the output is the converted XML file.
96 The overall processing pipeline of bioPDFX consists of the following five components, each
97 address an abovementioned challenge:

- 98 1. *Title and Abstract Correction*. We used the PMC Entrez e-utilities API to retrieve the
99 correct title and abstract from PubMed.
- 100 2. *Reference and Acknowledgment Correction*. To detect whether a paragraph is reference,
101 acknowledgement, or none of them, we designed two binary classifiers, one for reference
102 and the other for acknowledgement. For reference detection, we extracted features such as
103 density of year (e.g., “2016”), density of “et. al.”, and density of numbers followed by the
104 dot (e.g., “12.”). We normalized these features and train a classifier to predict whether a
105 paragraph is a reference text. For acknowledgment detection, we used a classifier with
106 features being TF-IDF followed by Latent Semantic Analysis (Deerwester et al. 1990) to
107 perform the detection.
- 108 3. *Special Character Correction*. The basic idea is to compare the results from PDFX and
109 that from Tesseract OCR to identify suspicious characters (i.e., different identified
110 characters within the same n-gram), and then apply Hidden Markov Model (HMM) (Baum
111 & Petrie 1966) with Viterbi inference algorithm (Viterbi 1967) and language model to
112 choose the “most probable” candidate character for the mismatch character by maximizing

113 the overall likelihood to recover the n-gram.

114 4. *Jumbling Error Correction*. We used the text extracted using Apache PDFBox to correct
115 the possible jumbling errors from PDFX. We first detected jumbled n-gram, and then
116 replaced them by the non-jumbled version from Apache PDFBox.

117 5. *Word Boundary Correction*. We adopted a simple dictionary lookup method by comparing
118 different parts (e.g., “can-” and “cer”) with the English dictionary in the Enchant
119 Spellchecking System (2015c) and merging them as required.

120 To evaluate the bioPDFX tool, we randomly extracted 100 biomedical literatures related to
121 Genome Wide Association Study (GWAS) (Hindorff et al. 2009; Welter et al. 2014) from PMC,
122 and used XML versions (i.e., NXMLs) of those 100 articles provided in PMC as gold standards.
123 We compare bioPDFX with the state-of-the-art PDFX tool (Constantin et al. 2013). As a use case,
124 we also compare the results of converting important information (p-value and number) and overall
125 text quality in the GWAS papers using bioPDFX and PDFX. We evaluate both individual
126 correction steps and overall conversion results. The experiment results show significant
127 improvement for all five types of correction tasks, especially for abstract,
128 reference/acknowledgement, special character, and jumbling error corrections.

129 Our contributions of this study are three-fold:

- 130 • We identified five issues (title/abstract, reference/acknowledgment, special character,
131 jumbling error, and word boundary) of the state-of-the-art PDF-to-XML tools as the very
132 first step of biomedical NLP pipeline.
- 133 • We integrated the output of four popular tools (PDFX, PMC Entrez e-utilities API,
134 Tesseract OCR, and Apache PDFBox) in bioPDFX and designed five corresponding

135 correction steps to solve the abovementioned issues.

- 136 • We evaluated bioPDFX and compare it to PDFX for each correction steps as well as a real-
137 world GWAS information extraction task, and the results demonstrated that bioPDFX can
138 improve all the five issues and increase the extraction accuracy for GWAS literatures.

139

140 **Related Work.** The rest of the section reviews the related studies along with the mention of tools
141 that we use in our approach. PDF text extraction involves extracting out all textual content from
142 the PDF file in an order that makes sense even without the formatting of the PDF. A variety of
143 tools are available that parse and convert PDF to text such as Apache PDFBox (2015b), Mozilla
144 PDF.js (2015d), and Adobe Acrobat SDK (2015a). Most text extractors use a page segmentation
145 algorithm to determine all the textual areas of the PDF based on width of the whitespace and from
146 these boundaries also determine a reading order. A major problem associated with the text based
147 extraction is because of variability of PDF publishing software. Some tools embed a special
148 character that is not included in a standard font by drawing a vector graphics over it. Thus, although
149 it might look perfectly correct to a reader, the information is not present in the text layer of the
150 PDF for text extractors to extract. In our bioPDFX tool, we integrated text extraction from Apache
151 PDFBox with other tools such as PDFX and optical character recognition tool Tesseract OCR to
152 overcome challenges like this.

153 XML transcription of a PDF file involves extraction of text, tables, images, etc. and then
154 tagging them with appropriate sections/regions to make an XML. The exact format of the
155 transcribed XML depends on the specific tool and might have a tag for title, abstract, introduction,
156 tables, etc. Some tools like PDFX (Constantin et al. 2013) use a schema similar to JATS/NLM
157 document-type definition (DTD) format which is a specific format for scientific articles and are

158 highly optimized around the organization of such PDF files. PDFX is a freely available rule-based
159 system which converts scholarly PDF articles to XMLs which have detailed sectional information
160 such as title, abstract, references, tables, and acknowledgment. It derives transcription parameters
161 from the relative font sizes, style and spacing of articles with respect to individual PDF files. In
162 our work, we integrate PDFX with other extraction tools to improve the extraction results for texts
163 like title, abstract, reference, and acknowledgement.

164 Text post-processing is a step that improves the correctness of the generated text. Several
165 tools, such as Tesseract OCR (Smith 2007), deploy techniques from Optical Character Recognition
166 (OCR) to increase the accuracy in the post-processing stage. Because most OCR engines classify
167 each character independent of other characters, it is important to post-process the text using some
168 contextual information such as statistical language models or syntactic and semantic rules. One
169 such method suggested by Zhuang and Zhu (Zhuang & Zhu 2005) reduces the candidate character
170 search space for characters by using the candidate distance information by an OCR engine in
171 conjunction with an n-gram based language model and a semantic lexicon. In another method
172 proposed by Velagapudi (Velagapudi 1999), contextual information was extracted by modeling
173 words/n-grams as sequences of characters and a Hidden Markov Model (HMM) was subsequently
174 used to decide the best sequence, using the given OCR output as the observation. Transition
175 probabilities of one character to another are calculated from a corpus and emission probabilities
176 are calculated from the output of the OCR engine on a corpus, to statistically record the error
177 patterns of the OCR engine. At the word level, Tong and Evans (Tong & Evans 1996) corrected
178 OCR text by modeling the text as sequence bigrams and then using HMM to compute the best
179 word sequence. The system learns the character level confusion probabilities for a specific OCR
180 engine and uses it to achieve better performance. In our study, we follow a similar approach to

181 post-process the output from PDFX using the Tesseract OCR engine.

182 In addition to above methodologies, several systems extract metadata such as authors, year of
183 publication, journal name, and bibliographical information from natural language text or PDFs
184 directly. CERMINE (Tkaczyk et al. 2014) uses supervised and unsupervised machine learning
185 techniques to extract parsed bibliographic references and metadata directly from PDFs. FLUX-
186 CiM (Cortez et al. 2007) uses unsupervised, non-template methods to extract citation components
187 from articles. Other PDF transcribers include LA-PDFText.(Ramakrishnan et al. 2012) Although
188 the abovementioned tools are available for PDF-to-text transcription, several issues are yet to be
189 dealt with, such as how to extract title and abstracts correctly, identify reference acknowledgement
190 sections in various format, recognize special characters correctly, detect and fix jumbling errors,
191 and correct word boundaries. PDFJailbreak (Garcia et al. 2013) is a communal project to create an
192 architecture and shared API to support open community development of semantic information
193 extractor from biomedical literature in PDF form.

194

195 2. Materials and Methods

196 The goal of this study is to extract high-quality XMLs from a biomedical literature PDFs to
197 make it better suited for our text mining tasks. Given a PDF file and its corresponding PubMed ID,
198 our system integrates the following four tools to generate output XML files:

199

- 200 • *PDFX* (Constantin et al. 2013) to transcribe XML of the given PDF
- 201 • *PMC Entrez e-utilities API* (Sayers et al. 2011) to retrieve text from PubMed
- 202 • *Tesseract OCR* (Smith 2007) to extract text extracted from the PDF file
- 203 • *Apache PDFBox* (2015b) to transcribe text of the given PDF

204

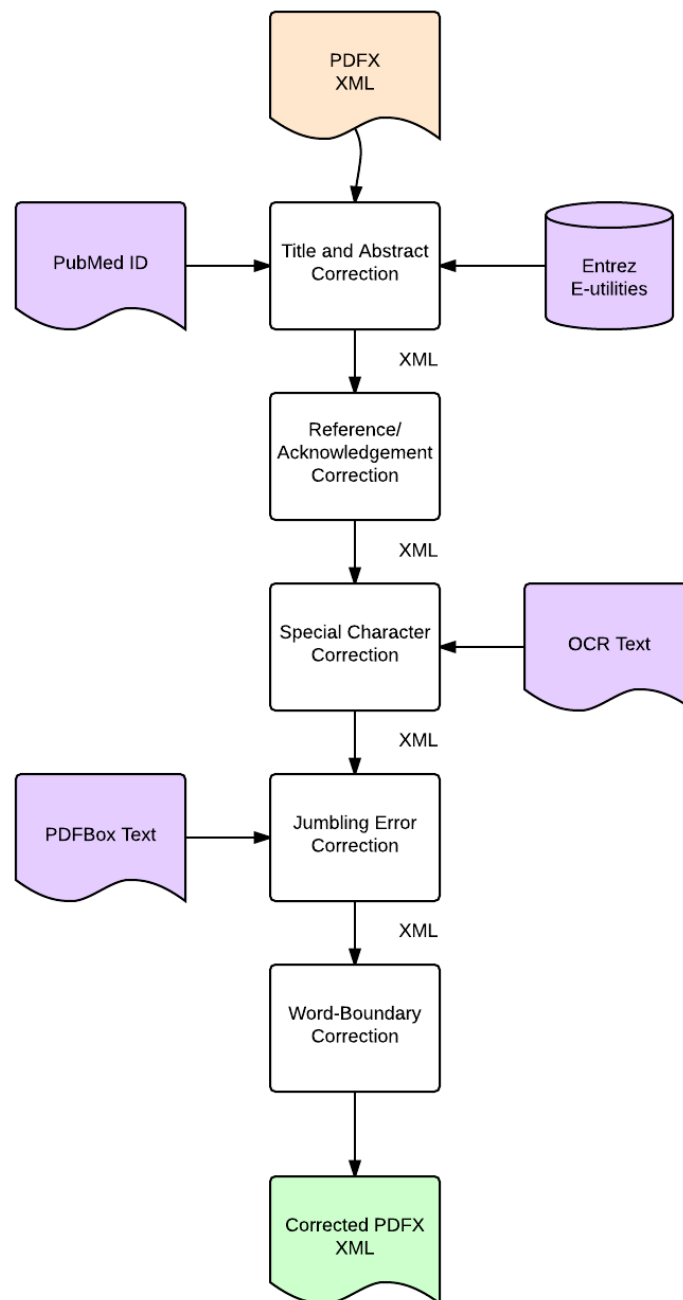
205 As shown in Figure 2, there are five correction stages (title/abstract,
206 reference/acknowledgment, special character, jumbling error, and word boundary corrections) in
207 the bioPDFX pipeline. Each stage is a filter and accepts an XML as input, giving a corrected XML
208 as output. We start from the XML generated by PDFX, correct title and abstract using PMC Entrez
209 e-utilities API, identify reference and acknowledgement texts, exploit Tesseract OCR to fix special
210 characters, utilize Apache PDFBox to correct jumbling errors, and recover word-boundaries as our
211 final stage. It should be noted that the pipeline design of bioPDFX allows different order of
212 correction stages, as well as adding new correction stages. The detail of each of the correction
213 stages is discussed in the following subsections.

214

215 **Figure 2. A flow diagram showing the stages of correction and their inputs.**

216 Stages are arranged in a pipelined fashion and take an XML as input and return a corrected XML

217 as output.



218

219

12

220 **2.1 Title and Abstract Correction**

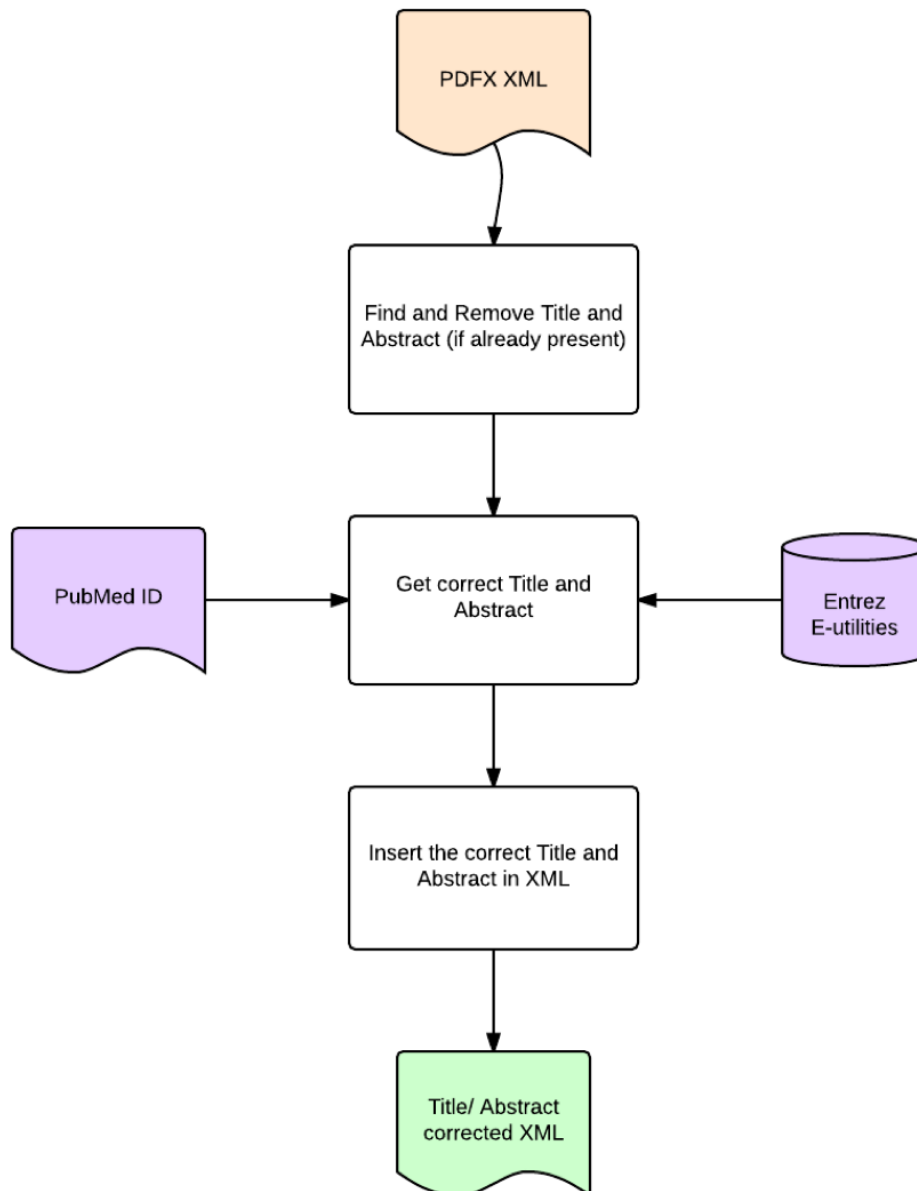
221 The flow diagram illustrating this step is given in Figure 3. As the input of this step, PDFX
222 performs well in detecting title and abstract for regular research articles; but in case of manuscripts
223 which might have contents other than title or abstract on their first page, PDFX struggles to identify
224 them correctly. To overcome this, we used the PMC Entrez e-utilities API (Sayers et al. 2011) to
225 retrieve the correct title and abstract (using PubMed ID of the paper) and substituted them for
226 whatever was present in PDFX XML. It should be noted that we applied PMC Entrez e-utilities
227 API because our focus is the biomedical literatures archived in PubMed Central.

228

229 **Figure 3. A flow diagram showing the title and abstract correction.**

230 The title and abstract correction pipeline, with PDFX XML as input and Title/Abstract corrected

231 XML as output.



232

233

234 2.2 Reference and Acknowledgement Correction

235 The format of references varies considerably among different publications ranging from a
236 separate section towards the end to individual references in footnotes continuing alongside the
237 content body. As PDFX is a purely rule-based system, it often falters in recognizing references
238 when the common pattern of having references in the end with a heading is not present. Similar
239 problems were noticed in identifying acknowledgements as well.

240 A key observation was that whenever PDFX identifies references/acknowledgement in a
241 paper, it is generally correct and correction was needed only in the case when it does not identify
242 either. We realized that the text content in references and acknowledgment showed clear patterns
243 in language, and because for our purpose, we did not need to extract individual bibliographic items
244 in the references section like PDFX does, classifying each section as a reference/acknowledgment
245 was enough. Also, it is important to note that the only requirement was to detect reference text,
246 which is much coarser grained and simpler than complete metadata extraction from journals, or
247 metadata extraction for each citation.

248 Although references and acknowledgment are usually clumped together, in this correction
249 step we designed two separate classifiers, as the patterns they show are different:

250 (1) For reference detection, we used PDFX XMLs which had a PDFX “reference” tag to train the
251 classifier (there are 451 such XMLs in our experiment), by taking reference text as positive
252 examples and rest of the text from the paper as negative examples. Some of the most prominent
253 features were density of year like numbers, density of “et. al” and similar strings, density of
254 numbers followed by the dot and so on. For each of the above features, the *ratio of the numbers*
255 *found versus total number of tokens in the text* was used instead of the raw numbers. We

256 normalized these features and train a Random Forest classifier (Breiman 2001) to predict
257 whether a paragraph is a reference text.

258 (2) For acknowledgment detection, we used PDFX XMLs which had a PDFX “acknowledgment”
259 tag to train the reference classifier (there are 451 such XMLs in our experiment). We exploited
260 another Random Forest classifier with features being TF-IDF followed by Latent Semantic
261 Analysis (Deerwester et al. 1990) (Truncated Singular Value Decomposition (Hansen 1987;
262 Kolda & O'leary 1998) in our experiment) to perform the detection. The features we adopted
263 are TF-IDF of the phrases such as “funded by” and “express gratitude”.

264

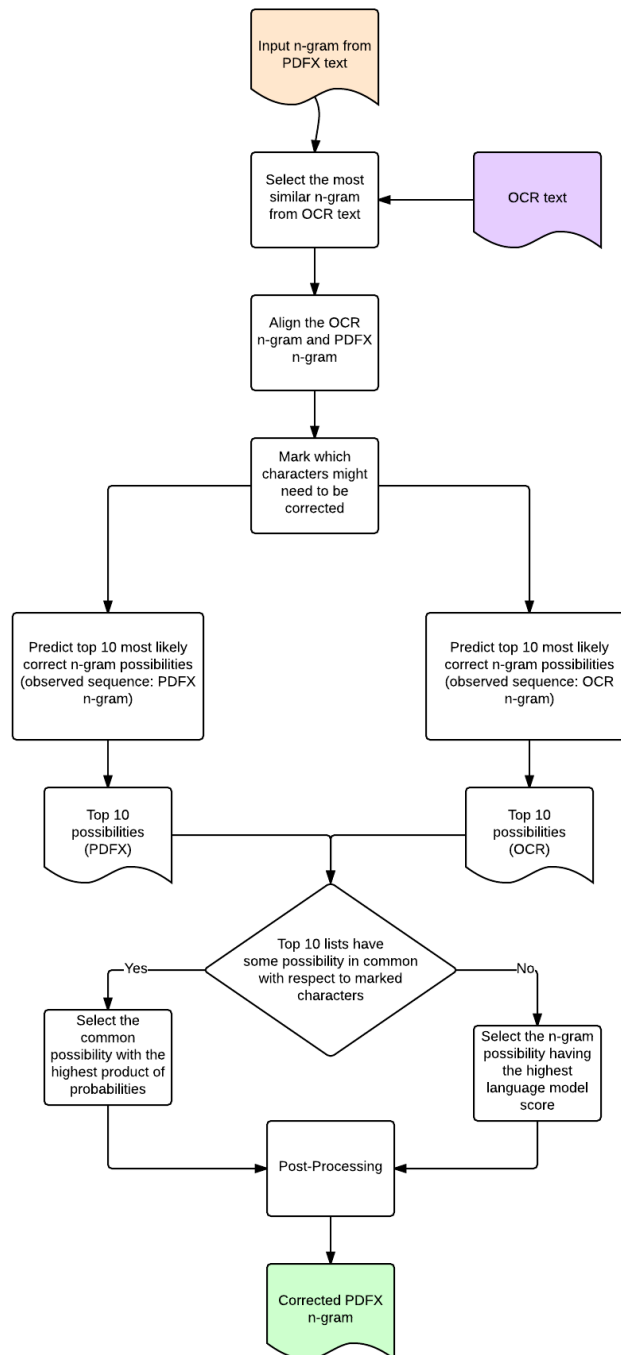
265 **2.3 Special Characters Correction**

266 PDFX and other PDF-to-text converters usually suffer from the problem of recognizing the
267 special characters (such as '!' and '@') depending on the publishing of the PDF files. The basic idea
268 for the special character correction step is that given a possibly erroneous n-gram from PDFX in
269 the token/word level, we recognized the corresponding n-gram in Tesseract OCR and then aligned
270 the two n-grams by their characters. Then, for differing characters we make a decision about what
271 character should be present in that position using Hidden Markov Model (HMM) (Baum & Petrie
272 1966) with Viterbi inference algorithm (Viterbi 1967) and language models. The flowchart for this
273 pipeline is summarized in Figure 4.

274

275 **Figure 4. A flow diagram for special character correction.**

276 This pipeline shows the steps taken in correcting special character errors in an n-gram.



277

278

17

279 2.3.1 Suspicious N-Gram and Character Identification

280 Before applying HMM, we first select similar character-level n-grams (we treat each character,
281 including the space, as unigram) from the output of OCR and PDFX, align the n-grams and then
282 mark suspicious characters in the n-grams. To select the most similar n-gram from OCR and PDFX
283 text, we used Levenshtein Distance (Levenshtein 1966) normalized for length to measure the
284 similarity between the two n-grams. The alignment of the two n-grams is done at the character
285 level and is based on a gene sequence alignment mechanism which is widely used (Hsu et al. 2008).
286 If the n-gram in OCR text is at least 0.6 (in Levenshtein Distance) as like the PDFX n-gram then
287 the characters that need to be replaced are marked. This is done by looking at the characters that
288 differ in the two n-grams after the alignment and selecting those characters which are not aligned
289 with whitespaces or digits. These *mismatch characters* are then fed to a HMM to predict top
290 possible correct character. The example in Figure 5 show the OCR, PDFX, and gold standard n-
291 grams. In this example, the mismatch characters in OCR n-gram are “α”, “u”, and “3”, while
292 that in PDFX n-gram are “a”, “a” and “s”.

293

294 **Figure 5. An example showing three types of n-gram: OCR, PDFX, and gold standard.**

295 The mismatch characters are shown in bold and underlined style.

1. OCR n-gram: “the **α**-val**ue** wa**3** less than 0.01”.
2. PDFX n-gram: “the **a**-val**ae** was less than 0.01”.
3. Gold standard n-gram: “the **p**-val**ue** was less than 0.01”.

296

297

298 **2.3.2 Prediction of Top-N Possible Correct N-Gram Lists for OCR and PDFX**

299 After marking the suspicious characters in the n-gram from OCR and the one from PDFX, we
300 predict the list of the most probable correct n-grams for OCR and the list for PDFX. Our solution
301 is to train an HMM (16) using a training corpus of aligned OCR, PDFX and gold standard texts,
302 and infer the “most probable” n-gram by the Viterbi algorithm.

303 While formulating this problem in terms of HMM Viterbi inference, one of our assumptions
304 is that characters in the words of meaningful sentences follow the Markov assumption to a good
305 extent, something which is already well proven by handwriting and speech recognition applications,
306 which have used HMMs to achieve good performance (10) (19). The other assumption is that
307 OCR and PDFX have a statistical pattern or bias in the character level errors that they make, and
308 that these patterns can be learned through a training corpus. The same technique of modeling
309 words as sequences of characters by HMM has been used to boost OCR accuracy in (24).

310 To formally state the problem, we first define a character x_t as the t-th letter in an OCR n-gram.
311 Therefore, each OCR n-gram can be represented as $X = x_1x_2\dots x_N$, which is a concatenation of N
312 characters. It should be noted that the character whitespace (i.e., “ ”) is also treated as a character.
313 In Figure 5, the number of characters in the OCR n-gram “the α -value was less than 0.01” is 30.

314 Next, we compute an *ambiguity dictionary*, which is a dictionary that records occurrences
315 of observed mismatched characters in the corpus. Specifically, the ambiguity dictionary is a list
316 of tuples of the form: (observed character, actual character, frequency). For example, a tuple (“ α ”,
317 “p”, 30) means that for 30 times when we saw “ α ” in OCR / PDFX text corpus, the actual
318 character in gold standard was “p”. Multiple tuples for an observed character, like “ α ” in our
319 example can be aggregated and written as:

320

321 “ α ”: [(“ p ”, 150), (“ α ”, 30), (“ a ”, 5)],

322

323 which means that for “ α ” in OCR / PDFX n-gram, the true character should be a “ p ” 150 times,
 324 “ α ” 30 times and “ a ”, 5 times.

325 Thus, each mismatch character (from OCR n-gram) x_t corresponds to a candidate-character set
 326 δ_t , which can be computed from the ambiguity dictionary. For non-mismatch characters, the
 327 candidate character set $\delta_t = \{x_t\}$, contains only one element which is exactly x_t . In our OCR n-
 328 gram example, the candidate-token set of the first character “ t ” is “ t ”, while that of the fifth
 329 character “ α ” may be {“ α ”, “ a ”, “ p ”, “ o ”, “ 0 ”}, for example. Therefore, our problem can be
 330 regarded as to choose the “most probable” candidate character for each mismatch character in
 331 OCR n-gram, which can maximize the overall likelihood to recover the gold standard n-gram.
 332 In other words, we want to find $Y^* = y_1 * y_2 * \dots * y_N = \operatorname{argmax}_Y P(Y | X)$, where $y_1 * y_2 * \dots * y_N$ are
 333 the “most probable” candidate characters given X , the pair of OCR and PDFX n-grams. We
 334 simplify the notation by using $P(Y)$ as the target probability that we want to maximize.

335 We apply Hidden Markov Model (HMM) to solve our problem. Based on the Markov
 336 assumption, the probability we want to maximize can be written as $P(Y) = P(y_1 y_2 \dots y_N) =$
 337 $P(y_1)P(y_2 | y_1)P(y_3 | y_2) \dots P(y_N | y_{N-1})$. The model of our HMM is defined as follows:

338

- 339 • Observation

340 $\sigma = p_1, p_2, \dots, p_M$, which is a set of all possible M values for all characters.

341

- State

343

$$A = \bigcup_{t=1}^N \delta_t = \{q_1, q_2, \dots, q_K\}$$

344

which is a set of all possible K values for all candidate characters.

345

- Initialization Probabilities

$\Pi = \{\pi_i \mid \pi_i = P(y_t = q_i)\}$, which is a set of initial probabilities for each state.

Therefore, $|\Pi| = K$. Suppose $L^1(q_i)$ is the conditional probability given by 1-gram

language model (trained from the gold standard n-grams) for each candidate-character

q_i (e.g., “a”), we compute π_i as follows:

351

$$\pi_i = L^1(q_i)$$

352

Note that

353

$$\sum_{i=1}^K \pi_i = 1$$

354

- Transition Probabilities

$A = \{a_{ij} \mid a_{ij} = P(y_{t+1} = q_j \mid y_t = q_i)\}$, which is a set of probabilities transiting from the

t-th candidate character y_t with state q_i , to the (t + 1)-th candidate-character y_{t+1} with

state q_j . Therefore, $|A| = K^2$. Suppose $L^2(q_i, q_j)$ is the conditional probability

359 given by 2-gram language model (trained from the gold standard n-grams) for two
 360 consecutive candidate-characters q_i and q_j (for example, “ α ”), we compute a_{ij} as
 361 follows:

$$362 \quad a_{ij} = L^2(q_i, q_j)$$

363

364 Note that

$$365 \quad \sum_{j=1}^K a_{ij} = 1 \text{ for all states } i = 1, 2, \dots, K$$

366

367 • Emission Probabilities

368 $B = \{b_{ik} \mid b_{ik} = P(x_t = p_k \mid y_t = q_i)\}$, which is a set of probability for the t-th
 369 candidate-character y_t with state q_i , to emit the t-th character x_t with observation
 370 p_k . Therefore, $|B| = K * M$. Suppose $C(p_k, q_i)$ is the count that a candidate-character q_i
 371 (e.g., “p” from gold standard n-gram) is recognized as a character p_k (e.g., “ α ” from
 372 OCR n-gram), and λ is a regularization constant, we compute b_{ik} as follows:

$$373 \quad b_{ik} = \frac{C(p_k, q_i) + \lambda}{\sum_{k'}^M (p_{k'}, q_i) + \lambda}$$

374 Note that, for all states $i = 1, 2, \dots, K$. In our example ambiguity dictionary for
 375 character “ α ” = [(“p”, 150), (“ α ”, 30), (“a”, 5)], the emission probability from “p” to
 376 “ α ” (i.e., the probability that we observe “ α ” in OCR n-gram, but the true character
 377 in gold standard n-gram is “p”) can be estimated as $((30 + \lambda) / (150 + 30 + 5 + \lambda))$.

378

379 Finally, we apply the Viterbi algorithm with our HMM model to decode the N most probable
380 Y^* . In this way, we have the top N most probable n-grams given an OCR n-gram as the
381 observation.

382 We follow the same process for PDFX as well, that is, treating the PDFX n-gram as the
383 observation and using OCR n-gram to see which characters are not in consensus. Similarly, the
384 ambiguity dictionary for PDFX was derived by applying the same method, except that we record
385 ambiguities for PDFX instead of OCR. Hence, we can compute another list of the top N most
386 probable PDFX n-grams along with their probabilities.

387

388 **2.3.3 Combination of Top-N Possible Correct N-Grams Lists**

389 Next, we select the most probable n-gram as the n-gram common in both lists (with
390 respect to marked characters only) which has the maximum combined probability, using the idea
391 of integrating gene mention tagging models (9). This is illustrated in Figure 6.

392

393 **Figure 6. The example showing two top-N list of the most probable NXML n-grams with**
 394 **their log probabilities given by Viterbi inference.**

395 Candidate number 5 in the left list (where PDFX n-gram is the observation) matches candidate
 396 number 0 in the right list (where OCR n-gram is the observation), which is the “chosen best”.

0. -97.107 (rs2252931 max P = 2.2610 29		-87.240 (rs2252931 max P=2.2×10 ⁻⁹
1. -97.410 (rs2252931 max P = 2.2×10 29		-90.024 (rs2252931 max P=2.2×10 ⁻⁹
2. -98.379 (rs2252931 max P = 2.2610 29		-90.306 (rs2252931 max P=2.2×10 ⁻⁹
3. -98.681 (rs2252931 max P = 2.2×10 29		-90.312 (rs2252931 max P=2.2×10 ⁻⁹
4. -98.901 (rs2252931 max P = 2.2610 -9		-90.507 (rs2252931 max P=2.2×10 ⁻⁹
5. -99.204 (rs2252931 max P = 2.2×10 -9		-90.974 (rs2252931 Max P=2.2×10 ⁻⁹
6. -99.725 (rs2252931 max P = 2>2610 29		-91.043 (rs2252931 max P=2.2×10 ⁻⁹
7. -99.741 (rs2252931 max P 9 2.2610 29		-91.903 (rs2252931 max P=2.2×10 ⁻⁹
8. -99.782 (rs2252931 Max P = 2.2610 29		-92.031 (rs2252931 Jax P=2.2×10 ⁻⁹
9. -99.814 (rs2252931 max P 5 2.2610 29		-92.061 (rs2252931 max P=2.2×10 ⁻⁹

***Chosen best:

397 0. -186.444400 -- (rs2252931 max P = 2.2×10⁻⁹

398

399 The PDFX n-gram in Figure 6 is “rs2252931 max P = 2.2610 29” and the OCR n-gram is:
400 “rs2252931 max P=2.2x10-9”. The characters highlighted in green/blue are the characters marked
401 for correction. We won’t apply correction to other characters, but instead just keep whatever we
402 observe in the PDFX n-gram because other mismatches are pairs of digits. The left list of 10
403 candidates is from treating the PDFX n-gram as the observation, while the right side shows the
404 candidates from treating the OCR n-gram as the observation. The numbers in negative are the log
405 probabilities for each sequence as given by the Viterbi algorithm.

406 The candidate numbered 5 in the left list matches the candidate numbered 0 in the right list
407 (for the highlighted characters), and this combination has the highest joint probability (although
408 other combinations also match like 5, 1 and 5, 4). Thus, we choose this candidate as the most
409 probable n-gram.

410 There can be a case where the two lists do not share a common candidate with respect to
411 marked characters. In that case, we use the n-gram language model trained on our corpus to
412 obtain the probabilities for each of these twenty candidates, and select the most probable n-gram
413 candidate. In case there is again a tie among all twenty, we take the top n-gram candidate from
414 the PDFX candidate list as the most probable n-gram. An example of this scenario is shown in
415 Figure 7, where the two lists do not share a candidate in common and the language model chooses
416 the best candidate in the OCR candidate list, as it scores the highest.

417 Finally, in the post-processing step, we pick the n-gram with the highest score, as the “most
418 probable” n-gram as the output of the special character correction step.

419

420 **Figure 7. The example showing two top-N list of the most probable n-grams with their log**
 421 **probabilities given by Viterbi inference.**

422 None of the pair of candidates in the lists match.

0.	-76.754	P = 9.0610	29	Beta 21.9		-75.916	P= 9.0×10 ⁻⁹	Beta -1.9
1.	-78.122	P = 9.0×10	29	Beta 21.9		-76.056	P= 9.0×10 ⁻⁹	Beta -1.9
2.	-78.548	P = 9.0610	-9	Beta 21.9		-76.194	P= 9.0×10 ⁻⁹	Beta -1.9
3.	-78.716	P = 9.0610	29	Beta -1.9		-77.299	P= 9.0X10 ⁻⁹	Beta -1.9
4.	-79.388	P 9 9.0610	29	Beta 21.9		-77.439	P= 9.0X10 ⁻⁹	Beta -1.9
5.	-79.460	P 5 9.0610	29	Beta 21.9		-77.577	P= 9.0X10 ⁻⁹	Beta -1.9
6.	-79.583	P 4 9.0610	29	Beta 21.9		-78.462	P= 9.0×10 ⁻⁹	Beta +1.9
7.	-79.710	P 3 9.0610	29	Beta 21.9		-78.582	P= 9.0-10 ⁻⁹	Beta -1.9
8.	-79.794	P 2 9.0610	29	Beta 21.9		-78.659	P= 9.0×10 ⁹	Beta -1.9
9.	-79.916	P = 9.0×10	-9	Beta 21.9		-78.723	P= 9.0-10 ⁻⁹	Beta -1.9

***Chosen best:

Nothing in common in the nbest lists...

Lang model score OCR top: -23.6266403198 P = 9.0×10⁻⁹ Beta -1.9

423 Lang model score PDFX top: -25.0832328796 P = 9.0610 29 Beta 21.9

424

425 **2.4 Jumbling Correction**

426 The correction of jumbling errors (the wrong order of text while converting PDF to text, as
427 illustrated in Figure 1) in the output XML from PDFX is done by using the textual output as
428 generated by Apache PDFBox. We start by making a non-overlapping n-gram set of PDFX XML
429 text and an overlapping n-gram set for Apache PDFBox text. Then, we use a simple exhaustive
430 method to check if a n-gram contains jumbling error. That is, if an overlapping n-gram is a
431 permutation of another n-gram, we consider that the n-gram contains jumbling error. Next, we
432 replace the suspicious n-gram by the non-jumbled version of n-gram from Apache PDFBox.
433 Finally, we string together the non-overlapping n-grams with a space to construct the jumbling
434 corrected text.

435

436 **2.5 Word Boundary Correction**

437 To correct word boundary error (e.g., “cancer” might become “can-” and “cer” if it is the last
438 word of a row, thus it should be combined and corrected as “cancer” instead of two words), we
439 adopted a simple dictionary lookup method by comparing different parts (e.g., “can-” and “cer”)
440 with the English dictionary in the Enchant Spellchecking System (2015c) and merging them as
441 required. Specifically, if the hyphenated word is not present in the English dictionary but the word
442 without the hyphen is present in the dictionary, we remove the hyphen from the word. If the word
443 with or without hyphen is not present in the dictionary, we keep it unchanged.

444

445

446 **2.6 Experiment Settings**

447 We conduct experiments to compare bioPDFX with state-of-the-art PDFX tool (Constantin et
448 al. 2013), and test the following two hypotheses:

449

450 (a) For each of the individual correction steps, bioPDFX provides better quality of PDF to
451 XML conversion.

452 (b) For overall pipeline, bioPDFX also provides better conversion quality.

453

454 In our experiment, we set $n=5$ for n-grams. For special character correction, we exploit
455 StochHMM (14) to calculate the top-N possible correct n-Grams lists by Viterbi inference, and
456 KenLM (8) to build language models. We set the regularization parameter to 10^{-6} .

457

458 **2.6.1 Dataset and Gold Standards**

459 To test the two hypotheses, National Human Genome Research Institute (NHGRI) provided
460 2,185 biomedical literatures (Jain et al. 2016) related to Genome Wide Association Study (GWAS)
461 (Hindorff et al. 2009; Welter et al. 2014). Among them, 941 are indexed by PubMed Central
462 (PMC), where both PDF and XML versions are available for download as we searched in 2015.
463 From these articles, we randomly sampled 100 pairs of PDF and XML versions (i.e., NXMLs)
464 from these 941 articles as gold standards for four of our correction tasks: title and abstract, special
465 characters, jumbling, and word boundary. For reference and acknowledgement, we built our own
466 gold standard from PDF files directly rather than using those present in NXMLs. The reason behind

467 this was incomplete reference sections in most of NXMLs and complete omission of
468 acknowledgement sections. For 70 randomly selected PubMed articles, we manually compiled
469 references and acknowledgement text from their PDFs to act as gold standard. For hypothesis (b),
470 we further manually labeled the p-value and numbers (which are import for GWAS literatures) as
471 our gold standard. We randomly split the data in to 50% training and 50% test.

472

473 **2.6.2 Evaluation Metrics**

474 We developed a scoring metric which measured how a corrected XML from our system
475 performs with respect to the corresponding NXML. In this metric, an n-gram based similarity
476 method is used which is common in plagiarism detection. The idea is to make two sets of n-grams
477 from NXMLs and from our system, and use them to measure the occurrence and correct order of
478 words. We then compute the F1-score of the two n-gram sets as our evaluation metrics for title
479 and abstract, reference and acknowledgement, special characters, and word boundary. We applied
480 macro F1-score for most of the correction tasks except special characters and word boundary (of
481 which micro F1-score is computed). This is because the number of special characters or words
482 with incorrect boundary in an article can vary considerably. For jumbling error, we cannot apply
483 the n-gram similarity methods directly, as the error may appear in the tables (as shown in Figure
484 1). Therefore, we compute the average counts of such errors over all the articles in the dataset as
485 our evaluation metrics. For hypothesis (b), we further compute the micro F1-score for the p-value
486 and numbers (for the same reason as special character and word boundary), and the macro F1-
487 score for the overall text quality (which is the n-gram similarity score for the whole document).

488

489 **3. Results**

490 The results for each of the five correction steps are shown in Table 1. That is, each step is
491 evaluated independently. As Table 1 depicts, in all correction steps bioPDFX outperforms PDFX,
492 especially for abstract, reference/acknowledgement, special character, and jumbling error
493 corrections.

494 On the other hand, Table 2 shows the scores at the end of the pipeline, where the corrections
495 are run in order of the pipeline as shown in Figure 2. Like the results of individual corrections
496 steps, bioPDFX in general perform better. The additional evaluation tasks (p-value, number and
497 overall text quality) are also included in Table 2. It should be noted that although no correction is
498 aimed at improving p-values, numbers or overall text quality, bioPDFX was still able to provide
499 higher conversion quality for these three tasks.

500

501 **Table 1. Stage wise quality score comparison of raw and corrected XMLs.**

Correction Step	PDFX	bioPDFX	Metric
Title	0.9135	0.9763	Macro F1
Abstract	0.5428	0.8920	Macro F1
Reference/Acknowledgment	0.6496	0.8026	Macro F1
Special Character	0.7071	0.8860	Micro F1
Jumbling	2.89	1.92	Average Error
Word-Boundary	0.7619	0.8053	Micro F1

502 Jumbling score is the average number of jumbling errors detected per paper (lower is better).

503

504 **Table 2. Final quality score comparison of raw and corrected XMLs at the end of the**
 505 **correction pipeline.**

Correction Step	PDFX	bioPDFX	Metric
Title	0.9135	0.9695	Macro F1
Abstract	0.5428	0.8877	Macro F1
Reference/Acknowledgment	0.6496	0.8026	Macro F1
Special Character	0.7071	0.8932	Micro F1
Jumbling	2.89	1.95	Average Error
Word-Boundary	0.7619	0.8189	Micro F1
P-value	0.5465	0.6615	Micro F1
Number	0.8733	0.8853	Micro F1
Overall Text Quality	0.9006	0.9107	Micro F1

506 Jumbling score is the average number of jumbling errors detected per paper (lower is better).

507

508 4. Discussion

509 Based on the experiment results, we show the experimental evidence to support the two
510 hypotheses. Next, we discuss our results for each correction step in detail:

511

512 • *Title and Abstract.* The correction scores for title and abstract is quite good in case of raw
513 XMLs and increases substantially in the final corrected XMLs. This is an expected
514 behavior as we are retrieving titles and abstracts using the PMC Entrez e-utilities API
515 directly which can be thought of as a gold standard itself. The F1-scores are still not perfect
516 due to syntactic differences between title/abstract text present in NXMLs and the
517 title/abstract text retrieved from e-utilities API. For example, the API does not seem to use
518 Unicode characters like “ \leq ”, “ \geq ” and “ \sim ” and would instead have the phrases “less than or
519 equal”, “greater than or equal” and “approximately” in their place. Furthermore, there
520 were cases where the data from NXMLs was incorrect because of errors from
521 authors of those papers like missing title/abstracts and other similar inconsistencies
522 in punctuation. Apart from this, we can consider the titles and abstracts to be near perfect in
523 the corrected XMLs (assuming e-utilities API gives correct results). Some other errors
524 are introduced due to accumulation of errors from previous stages of corrections in the
525 pipeline, but as can be seen by comparing stage-wise and end of pipeline scores in Table 1
526 and Table 2, they are quite small and can be neglected.

527 • *Reference and Acknowledgement.* Reference/Acknowledgment scores also improved
528 considerably and lead to better scores for the remaining stages of the pipeline, by removing
529 a significant number of the False Positives, as can be observed by comparing the stage-

530 wise-results in Table 1 and final-results in Table 2.

531 • *Special Character*. We could improve special character F1-score by a considerable margin
532 through our technique of combining HMMs, n-gram alignment, OCR, and language
533 models.

534 • *Jumbling Error*. The average counts of jumbling errors are fewer for bioPDFX then for
535 PDFX. It should be noted that there is no guarantee that jumbling errors occur only over 5-
536 grams, or equivalently within 5 words. In the same document, there might be jumbling
537 errors spanning up to arbitrarily large and varied n. A better approach might be to start with
538 a much larger n and attempt to detect jumbling from there and progressively reduce n one
539 by one, however this approach is computationally intensive. Even operating on 5-grams
540 for jumbling error correction step takes a significant amount of time compared to other
541 correction steps.

542 • *Word Boundary*. The F1-score of word boundary correction using bioPDFX are also
543 improved modestly comparing to PDFX. We also attempted to use a language model (along
544 with the English dictionary) to increase the range of corrections possible for biomedical
545 terms, but it did not give us any performance improvement. We believe this is because
546 word boundary errors are comparatively infrequent compared to total number of words in
547 an article, and because most words in an article are English words or numbers, most word
548 boundary errors happen in English words, which are already corrected by the English
549 dictionary.

550 Also, the usability of bioPDFX can be largely increased by a multiple-files uploading
551 functionality, so that users can submit many PDF files of interest for conversion instead of

552 uploading one-by-one. We are currently implementing an interface for such a multiple-uploading
553 feature (Figure 8) for a biocuration tool that will include this feature to maximize the utility of
554 bioPDFX.

555

556

557 **Figure 8. The interface for multiple-files uploading.**

558 The usability of bioPDFX can be largely increased by integrating it within a biocuration tool.

The screenshot displays the 'Stage One : View Files' interface within the 'GWAS CURATION' system. The user is identified as 'Chun Nan' and is logged in. The interface features a sidebar with navigation options: DASHBOARD, NEW PROJECT, SETTINGS, and STATISTICS. The main content area includes a search bar and a table of files. Below the table, there is a progress indicator showing '5 rows visible' and a pagination control. A preview area displays four selected PDF files, each with a 100% progress bar and a checkmark. The bottom bar indicates '4 files selected' and provides buttons for 'Remove Files', 'Upload Files', 'Browse Files', and 'Proceed to Stage Three'.

FILE ID	NAME	TYPE	SIZE	ACTIONS
<input type="checkbox"/> 364c27a1cb3783c7fbd2abe6d193e8b2	19122664.xml	text/xml	87231	✕
<input type="checkbox"/> 2561e2c428a6003bfb89c9117e8202c0	19597492.xml	text/xml	69648	✕
<input type="checkbox"/> 34e5728bd9ce0df5811163b21e15bb88	19412176.xml	text/xml	80656	✕
<input type="checkbox"/> 14fceaed83649712c6066a354488d39e	19915575.xml	text/xml	90928	✕
<input type="checkbox"/> fd98d38c6292030b555cd7594f3314dc	19448621.xml	text/xml	89376	✕

5 rows visible

4 files selected

Remove Files Upload Files Browse Files

Proceed to Stage Three

559

560

561 **5. Conclusions**

562 We have developed a tool to convert PDF files into XML format specifically for biomedical
563 articles, targeting the five challenges of state-of-the-art converters: title/abstract,
564 reference/acknowledgement, special character, jumbling error, and word boundary. Overall, the
565 whole pipeline developed in this paper makes the published literature in form of PDF files much
566 better suited for text mining tasks while slightly improving the overall text quality as well.
567 Although the tool discussed in this paper is trained to be specifically optimized for our target
568 corpus (GWAS), we believe that these techniques can be applied to other kinds of literatures as
569 well.

570 It should be noted that although we use PDFX XMLs as the main text sources, and use Apache
571 PDFBox and OCR text as auxiliary text sources, bioPDFX is general and is not dependent on these
572 text sources. Given a primary text source to correct substitution errors, other secondary text sources
573 can also be exploited in the correction steps.

574 There are some limitations of the approach, but we believe that these can be easily overcome
575 in different domains and by investigation of several different text sources. These limitations are
576 summarized as below:

577

- 578 • Title and abstract corrections depend on the PMC Entrez e-Utilities API
- 579 • As we model n-grams as sequences of characters and due to the Markov assumption, we
580 cannot correct errors in which one character is replaced by multiple ones, or is deleted
- 581 • Jumbling error correction is only an approximation since we fix $n=5$

582

583 In the future, we plan to run each of the text sources in parallel further improving the overall
584 performance, for which our initial result was reported in (Goyal et al. 2016). Also, we plan to
585 explore more challenging issues in the PDF to XML conversion process. Finally, we plan to extend
586 our experiments on more diverse types of large-scale corpuses.

587

588 **Acknowledgements**

589 We would like to thank Dr. Lucia Hindroff of NHGRI for her generous support of this project by
590 providing us PDF copies of test articles and the curation guidelines of the Catalog of GWAS to
591 make this research possible. We also would like to thank Dr. Helen Parkinson, Dr. Jackie
592 MacArthur and their team members at EBI for their assistance.

593

594 **References**

- 595 2015a. Adobe Acrobat SDK. Available at <http://www.adobe.com/devnet/acrobat/overview.html>.
- 596 2015b. Apache PDFBox. Available at <https://pdfbox.apache.org/index.html>.
- 597 2015c. Enchant Spellchecking System.
- 598 2015d. Mozilla PDF.js. Available at <https://mozilla.github.io/pdf.js/>.
- 599 2015e. National Center for Biotechnology Information (NCBI), National Library of Medicine (NLM),
600 PubMed Central (PMC) Open Access Subset. Available at <http://www.ncbi.nlm.nih.gov/pmc/tools/ftp/>.
- 601 Baum LE, and Petrie T. 1966. Statistical inference for probabilistic functions of finite state Markov chains.
602 *The annals of mathematical statistics* 37:1554-1563.
- 603 Berg ØR. 2011. High precision text extraction from PDF documents.
- 604 Breiman L. 2001. Random forests. *Machine Learning* 45:5-32.
- 605 Constantin A, Pettifer S, and Voronkov A. 2013. PDFX: fully-automated PDF-to-XML conversion of
606 scientific literature. Proceedings of the 2013 ACM symposium on Document engineering. Florence,
607 Italy: ACM. p 177-180.
- 608 Cortez E, da Silva AS, Gonçalves MA, Mesquita F, and de Moura ES. 2007. FLUX-CIM: flexible
609 unsupervised extraction of citation metadata. Proceedings of the 7th ACM/IEEE-CS joint
610 conference on Digital libraries: ACM. p 215-224.
- 611 Deerwester S, Dumais ST, Furnas GW, Landauer TK, and Harshman R. 1990. Indexing by latent semantic
612 analysis. *Journal of the American Society for Information Science* 41:391-407. 10.1002/(sici)1097-
613 4571(199009)41:6<391::aid-asi1>3.0.co;2-9
- 614 Garcia A, Murray-Rust P, Burns G, Stevens R, Tkaczyk D, McLaughlin C, Belin A, Di Iorio A, Garcia L,
615 and Gruson-Daniel C. 2013. PDFJailbreak-a communal architecture for making biomedical PDFs
616 semantic. *Proceedings of BioLINK SIG 2013*:63.
- 617 Goyal A, Singh A, Bhargava S, Crawl D, Altintas I, and Hsu C-N. 2016. Natural Language Processing
618 using Kepler Workflow System: First Steps. *Procedia Computer Science* 80:712-721.
- 619 Hansen PC. 1987. The truncatedsvd as a method for regularization. *BIT Numerical Mathematics* 27:534-
620 553.
- 621 Hindorff L, Sethupathy P, Junkins H, Ramos E, Mehta J, Collins F, and Manolio T. 2009. Potential etiologic
622 and functional implications of genome-wide association loci for human diseases and traits.
623 *Proceedings of the National Academy of Sciences* 106:9362-9367. citeulike-article-id:4806365
624 doi: 10.1073/pnas.0903103106
- 625 Hsu C-N, Chang Y-M, Kuo C-J, Lin Y-S, Huang H-S, and Chung IF. 2008. Integrating high dimensional
626 bi-directional parsing models for gene mention tagging. *Bioinformatics* 24:i286-i294. citeulike-
627 article-id:12926636
- 628 Jain S, Tumkur K, Kuo T-T, Bhargava S, Lin G, and Hsu C-N. 2016. Weakly Supervised Learning of
629 Biomedical Information Extraction from Curated Data. *BMC Bioinformatics*.
- 630 Kolda TG, and O'leary DP. 1998. A semidiscrete matrix decomposition for latent semantic indexing
631 information retrieval. *ACM Transactions on Information Systems (TOIS)* 16:322-346.
- 632 Levenshtein VI. 1966. Binary codes capable of correcting deletions, insertions and reversals. Soviet physics
633 doklady. p 707.
- 634 Ramakrishnan C, Patnia A, Hovy E, and Burns GA. 2012. Layout-aware text extraction from full-text PDF
635 of scientific articles. *Source code for biology and medicine* 7:1.
- 636 Sayers EW, Barrett T, Benson DA, Bolton E, Bryant SH, Canese K, Chetvernin V, Church DM, DiCuccio
637 M, and Federhen S. 2011. Database resources of the national center for biotechnology information.
638 *Nucleic acids research* 39:D38-D51.
- 639 Smith R. 2007. An overview of the Tesseract OCR engine.
- 640 Tkaczyk D, Szostek P, Dendek PJ, Fedoryszak M, and Bolikowski L. 2014. CERMINE--Automatic
641 Extraction of Metadata and References from Scientific Literature. Document Analysis Systems

- 642 (DAS), 2014 11th IAPR International Workshop on: IEEE. p 217-221.
- 643 Tong X, and Evans DA. 1996. A statistical approach to automatic OCR error correction in context.
- 644 Proceedings of the fourth workshop on very large corpora. p 88-100.
- 645 Velagapudi P. 1999. Using hmms to boost accuracy in optical character recognition. Proceedings of SPIE,
- 646 27th AIPR Workshop: Advances in Computer-Assisted Recognition: Citeseer. p 96-104.
- 647 Viterbi A. 1967. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm.
- 648 *IEEE TRANSACTIONS ON INFORMATION THEORY* 13:260-269.
- 649 Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, Klemm A, Flicek P, Manolio T, Hindorff
- 650 L, and Parkinson H. 2014. The NHGRI GWAS Catalog, a curated resource of SNP-trait
- 651 associations. *Nucleic acids research* 42:D1001-D1006. citeulike-article-id:12826571
- 652 doi: 10.1093/nar/gkt1229
- 653 Zhuang L, and Zhu X. 2005. An OCR post-processing approach based on multi-knowledge. International
- 654 Conference on Knowledge-Based and Intelligent Information and Engineering Systems: Springer.
- 655 p 346-352.
- 656