

# CHECKBOX DETCTION

## Abstract

The Portable Document Format (PDF) is the most used file format for online documents. The absence of effective means to extract text from these PDF files in a structured manner presents a significant challenge for developers. Our paper describes the construction and performance of a system that extracts text blocks from PDF-formatted documents and classifies them into logical units based on prior hand-made annotations that characterize specific sections. We also introduce new state-of-art check-box detection schema using RetinaNet deep-learning neural network.

## Problem Statement

The purpose of the document analysis project is to create a web application, which will disaggregate multiple pdf format documents of a same type and reaggregate the information within in a manner that would make it identifiable, **the** retrievable and easily readable. Generally, data in **the** pdf documents is semi-structured. As opposed to well-structured data, which conforms to a schema or data model and can be queried using structured query language to answer questions, it does not adhere to any rigorous format.

## Solution

To receive structured data as described in problem statement, we propose the following steps:

<b>Step 1:</b> Annotate single document of <b>an</b> each type	<b>Step 2:</b> Extract data from multiple documents in accordance with initial annotation	<b>Step 3:</b> Identify only information from checked boxes for extraction	<b>Step 4:</b> Create structured data in a csv format
---	--	---	--

**Step 1** of the algorithm is to annotate a single document of a given type. See **Figure 1** for examples of uploaded and annotated documents.

ADVERSE DRUG REACTION REPORTING FORM	
Sr. No	REPORT ON SUSPECTED SERIOUS ADVERSE DRUG REACTION
<div>For Report to Drugs Controller Pak Secretariat, Block C, Ministry of Health,</div>	
<b>1. PARTICULARS OF PATIENT</b>	
Name of patient <u>Rose Matthew</u>	
Age <u>53</u>	Weight (kg) <u>UKWN</u> Patient address <u>Ohio</u>
Sex <input type="checkbox"/> Male <input checked="" type="checkbox"/> Female	Race <u>Not Known</u>
Pregnant <input type="checkbox"/> Yes <input type="checkbox"/> No <input checked="" type="checkbox"/> Not applicable	
Relevant Medical History <u>Patient stated she had "normal" stiffness in left deltoid muscle day after vaccine, but on 12/31 (3 full weeks after the vaccine), stated still having tightness in arm and throbbing when she lays down, intermittently.</u>	
<b>2. ADVERSE EVENT</b>	
Reason for reporting	
<input type="checkbox"/> Requires or prolongs hospitalization	<input type="checkbox"/> Life threatening <input type="checkbox"/> Death
<input type="checkbox"/> Permanently disabling or incapacitating	<input type="checkbox"/> Congenital anomaly <input type="checkbox"/> Overdose
<input type="checkbox"/> Other (Please Specify) <u>Injection site reaction</u>	



**ADVERSE DRUG REACTION REPORTING FORM**

**REPORT ON SUSPECTED SERIOUS ADVERSE DRUG REACTION**

For Report to  
Drugs Controller  
Pak Secretariat, Block C,  
Ministry of Health,  
...

**1. PARTICULARS OF PATIENT**

Sr. No. [676023]

Name of patient [Emily Whitter]

Age [65 Yrs] Weight (kg) [150 lbs] Patient address [UT]

Sex ☐ Male ☒ Female

Race [White] Caucassian

Pregnant ☐ Yes ☒ No ☐ Not applicable

Relevant Medical History [Gout, High Cholesterol, High A1C; Medicine allergies include: Allopurinol, DARVON, Levofloxacin, Meloxicam, Omeprazole, Penicillins, Penicillamine, PHENERGAN]

**2. ADVERSE EVENT**

Reason for reporting

☐ Requires or prolongs hospitalization ☐ Life threatening ☐ Death

☐ Permanently disabling or incapacitating ☐ Congenital anomaly ☐ Overdose

☒ Other (Please Specify) See next page for details

**Figure 1:** An example of an unannotated document (upper left) and annotated document (lower right).

The document should be annotated so that the key and value pairs would be identifiable, I.e. Name: Smith, Age: 58 years, where key is the question, and value is the corresponding answer. In picture 1 it is apparent that we have chosen to annotate keys with red boxes, while the values with blue boxes.

**Step 2.** Using a python library fitz, we can detect the red and blue boxes and extract their corresponding the coordinates. We can also detect the text within the boxes and extract them into any python data structure. Thus, keeping the coordinates of the boxes, we can receive the information from the boxes from any document of the same type. As you can see in Figure 3, the data from three analogous documents for three different patients was received having coordinates of annotation boxes of a single document.

	Sr. No	Name of patient.	Age	Weight (kg)	Patient address	Sex	Race	Pregnant	Relevant Medical History	Reason for reporting	...
0		Rose Matt h...	53	U K W N	Ohio	Male Fem a l ...	Not Kno wn	Yes No Not...	Patient st...	Requires o...	...
1	[676023]	[Emily Whitter]	[65 Yrs]	[150 lbs]	[UT]	[Male, Female]	[White Caucassian]	[Yes No Not applicable]	[Gout, High Cholesterol, High A1C; Medicine al...	[Requires or prolongs hospitalization Life thr...	...
2	[676158]	[Hayden Blank]	[59]	[U]	[New York City]	[Male, Female]	[African American]	[Yes No Not applicable]	[Season Allergies (Pollen), Other Medication ...	[Requires or prolongs hospitalization Life thr...	...

**Figure 3:** Structured data received from three documents of the same type.

However, to retrieve only information from a checked box is yet another challenge. In figure 3 we can observe for section “Sex” both options – male and female. Information from the annotated boxes is retrieved fully irrespective of the checked box in front of only one of the options.

**Step 3.** To detect which box is checked, we suggest using RetinaNet deep learning network. We propose a completely new solution to the problem of checkbox detection: instead of detecting the box itself, we trained the network in a manner, that it could identify the bounding box of the text near the checked box (see Figure 4). This facilitates the work, because it reduces the number of steps in the code. It also gives a possibly higher accuracy, since we don’t have to rely on engineering solutions to find out whether the check box is on the left of the text or on the right.

☒ Male  
☐ Female

Life threatening <input type="checkbox"/>	Caused or prolonged hospitalisation <input type="checkbox"/>	Caused disability or incapacity <input checked="" type="checkbox"/>	Caused birth defect <input type="checkbox"/>	*N/A (not serious) <input type="checkbox"/>
---	--	---	--	---

Figure 4. Detection of the text near the checked box using RetinaNet deep learning network.

We have also considered cases where there is an additional information not only near the checkboxes, or the words selected are circled instead (Figure 5). For instance, in the Figure 5, “56 years” is the output of our program.

of 1,000  
 mins/ hours/ days/ months/ years  
 (please circle)

Figure 5. Detection of circled words and additional words.

After coordinate detection of the bounding box of the text, we proceed to the **Step 4**. The image is cropped according to coordinates and the image is passed to tesseract OCR. OCR transferred text is then inserted to the corresponding space in the data structure shown in Figure 3.

## Data Preprocessing for RetinaNet Neural Network

The data for the network was collected from various documents containing different type of checkboxes (Figures 4, 5). Checkmarks inside the boxes also vary, being X-shaped or V-shaped. Some noises were added to the pictures like dots of different colors to avoid overfitting. Negative examples included both – figures with all unchecked boxes and figures without any checkboxes on them.

For positive data a code was written using pytesseract python library to give as an output an hocr file. Hocr is an open standard of data representation for formatted text obtained from, optical character recognition, which among many other things contains the bounding box coordinates of the words. For each figure a bounding box coordinates of the words were found and were given as labels of the data. Three different type of classes were created: one for words near the checkboxes for cases represented in picture 4, and two additional classes for cases presented in Figure 5: one for the text input without checkboxes, and the other for circled words. We named the classes yes, ok and ellipse. You can see an example of preprocessed data in Figure 6.

Checkbox image	X <sub>min</sub>	Y <sub>min</sub>	X <sub>max</sub>	Y <sub>max</sub>	Class name
Mild <input type="checkbox"/> Moderate <input checked="" type="checkbox"/> Severe <input type="checkbox"/>	505	18	670	55	tick
<div>95 mins/ hours/ days/ months/ <u>years</u> (please circle)</div>	80	25	135	65	notation
	660	10	77	65	circle
Certain <input type="checkbox"/> Probable <input checked="" type="checkbox"/> Possible <input type="checkbox"/> Unlikely <input type="checkbox"/> Unclassifiable <input type="checkbox"/>	330	15	480	50	tick
Certain <input type="checkbox"/> Probable <input type="checkbox"/> Possible <input type="checkbox"/> Unlikely <input type="checkbox"/> Unclassifiable <input checked="" type="checkbox"/>	1340	15	1580	50	tick
Life threatening <input type="checkbox"/> Caused or prolonged hospitalisation <input type="checkbox"/> Caused disability or incapacity <input type="checkbox"/> Caused birth defect <input checked="" type="checkbox"/> *N/A (not serious) <input type="checkbox"/>	1385	5	1620	95	tick
<div>53 mins/ hours/ days/ months/ <u>years</u> (please circle)</div>	80	25	135	65	notation
	660	10	770	65	circle
Recovered fully <input type="checkbox"/> Recovering <input type="checkbox"/> Not recovered <input checked="" type="checkbox"/> Unknown <input type="checkbox"/> Fatal <input type="checkbox"/> Date & Cause of death:.....	700	10	890	95	tick
<div>19 mins/ hours/ days/ months/ <u>years</u> (please circle)</div>	85	25	140	65	notation
	660	10	770	65	circle
<div>55 mins/ hours/ days/ months/ <u>years</u> (please circle)</div>	70	25	125	65	notation
	660	10	770	65	circle
Yes <input type="checkbox"/> No <input type="checkbox"/> Unknown <input type="checkbox"/> *N/A (drug continued) <input checked="" type="checkbox"/>	840	25	1190	70	tick
Mild <input type="checkbox"/> Moderate <input type="checkbox"/> Severe <input checked="" type="checkbox"/>	1100	20	1225	55	tick

**Table 6.** Preprocessed data for RetinaNet deep learning neural network.

Each of the data lines in **Table 6** contains of path to image where it is located, x,y coordinates of upper left and lower right corners of the bounding boxes of the label data (cells 2 through 5), and their respective classes (cell 6).

The data was checked for avoiding possible mistakes and making it work for RetinaNet<sup>1</sup> through keras-retinanet/keras\_retinanet/bin/debug.py script designed for debugging images by RetinaNet model developers. To make text detectable by the neural network anchor parameters of anchor boxes were set to ratios = np.array([0.01,0.03,0.09,0.27, 1], keras.backend.floatx()), scales = np.array([2<sup>0</sup>, 2<sup>(1.0 / 3.0)</sup>, 2<sup>(2.0 / 3.0)</sup>, 3.5] in keras-retinanet/keras\_retinanet/utils/anchors.py<sup>1</sup> python script. All the coordinates were verified through debug.py

## RetinaNet Training

One of peculiarities of RetinaNet<sup>2</sup> model is the usage of Focal Loss ( $FL(p_t) = -(1 - p_t)^{\gamma} * \log(p_t)$ ) function during the classification. This function increases the overall contribution of positive examples in the training process. The Focal Loss is designed to address the one-stage object

<sup>1</sup> <https://github.com/fizyr/keras-retinanet>

<sup>2</sup> <https://arxiv.org/pdf/1708.02002.pdf>

detection scenario in which there is an extreme imbalance between foreground and background classes during training (e.g., 1:1000).

Training data contained 1205 samples. The parameters of training are as follows:

batch\_size= 1, steps= 1249, epochs= 50. We used pretrained COCO weights for RetinaNet<sup>1</sup>.

We used random-transform and no-resize parameters while training<sup>1</sup>. Random transform usage allows not to use min\_rotation, max\_rotation, min\_translation, max\_translation, min\_shear, max\_shear, min\_scaling, max\_scaling, flip\_x\_chance, flip\_y\_chance transformation parameters on initial pictures. No-resize is designed to feed the network with pictures of original sizes, and not set one size for all of them

## **Conclusion**

RetinaNet model was used for detecting/predicting checkboxes and their corresponding values. RetinaNet outputs the probabilities of checkboxes (i.e. the probability of existence of a checkbox), and we consider boxes as checked if the probabilities are higher than 0.5. It has been shown that our model detects checkboxes with high accuracy, approximately 90% on a given test data.

05 March 2020