

# Development of a Python Parser to Extract Information from a Scientific Paper in PDF format and write an output file

Clément Le Bihannic-Bertôt  
*Université Bretagne Sud 56000 Vannes, France*  
`le-bihannic-bertot.e2000371@etud.univ-ubs.fr`

Clément Cano  
*Université Bretagne Sud 56000 Vannes, France*  
`cano.e2000507@etud.univ-ubs.fr`

Théo Genette  
*Université Bretagne Sud 56000 Vannes, France*  
`genette.e2001218@etud.univ-ubs.fr`

Evann Perrinel  
*Université Bretagne Sud 56000 Vannes, France*  
`perrinel.e2000399@etud.univ-ubs.fr`

Iwan Fouquet  
*Université Bretagne Sud 56000 Vannes, France*  
`fouquet.e2002491@etud.univ-ubs.fr`

May 2023

## Abstract

L'IRISA avait besoin d'un analyseur de documents texte générés à partir d'un PDF. L'objectif était de découper le mieux possible les sections d'un article scientifique afin que les chercheurs gagnent du temps dans la lecture des articles. L'IRISA souhaitait un système écrit dans un langage de programmation adéquat et fonctionnant sur un système GNU/Linux. Pour répondre à ce besoin nous avons formé une équipe de 5 étudiants et travaillé en mode agile avec la méthode SCRUM.

## 1 Introduction

Notre objectif était de convertir un fichier PDF vers un fichier au format XML ou TXT. Pour cela, nous avons décidé que la manière la plus simple et la plus rapide de mettre en place notre projet était

d'utiliser le langage Python. À la fois pour sa flexibilité et pour l'intégration très simple des expressions régulières dans ce langage, élément central dans notre projet. En plus d'une version utilisable via un terminale, nous disposons d'une interface graphique bien plus simple à utiliser et toujours réalisée en Python à l'aide du module TKinter et CustomTKinter. TKinter est la librairie graphique de base du langage Python, elle ne permet de faire que des interfaces très simples mais c'était pour nous largement suffisante puisque nous devions sélectionner un fichier sur le système, choisir entre deux radio boutons, et sauvegarder un fichier. CustomTKinter nous a servi à rendre plus agréable visuellement l'interface créée à l'aide de TKinter. Avec cette librairie il a été par exemple possible de rendre l'interface plus sombre afin de surligner en couleurs les aspects importants sans pour autant perdre en esthétique, et de manière générale, de produire une interface moderne et lisible.

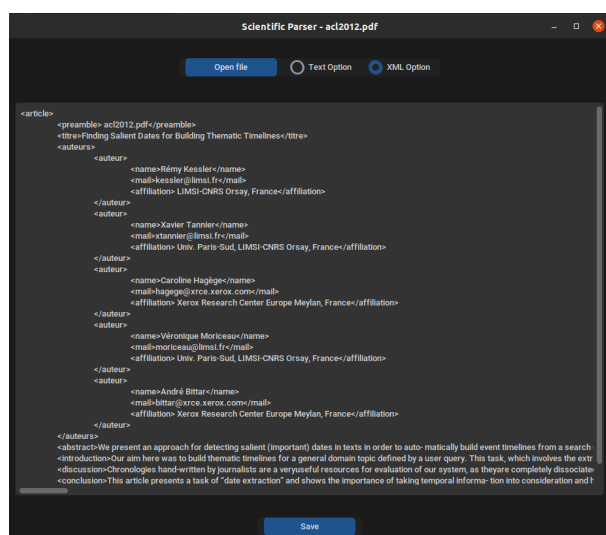


Figure 1: Interface Graphique de l'application

## 2 Explication du système

Concernant le code écrit dans ce projet, de manière générale, nous avons élaboré une expression régulière pour chaque partie présente dans tout article scientifique (abstract, results, conclusion, ...). Cette dernière était plus ou moins complexe selon les parties afin de capturer le texte concerné. Prenons par exemple la partie Abstract. Nous avons tout d'abord écrit une expression régulière comportant des groupes nommés

pour retrouver simplement l'information qui nous intéressait mais aussi de vérifier la présence de certains éléments comme un point à la suite du numéro de la section. Pour simplifier le debug de notre expression régulière nous avons ajouté dans chaque méthode une boucle for permettant d'afficher chaque groupe capturé et le texte associé. Enfin, nous avons récupéré les index des groupes qui nous intéressaient (souvent le texte uniquement) et nous avons appliqué de nouveau l'expression régulière sur tout le texte.

```

1 def getAbstract(text, file_name=""):
2     """
3     Extracts the abstract from the PDF using multiple regex to achieve the highest accuracy
4     possible
5     """
6     if(file_name=="IPM1481.pdf"):
7         regex=r"(?:([1-9]+?.?)|([IVX]*.))?(\\s+)?(?: (Abstract)|(ABSTRACT)|abstract)((\\s+)? *\\
8         n?)(?P<text>(?:.\\n)*?)^(?([1-9]+?.?)|([IVX]*.))?.s+(Introductio.n|INTRODUCTION|
9         A r t i c i a l ))| ntroductio .n\\n)"
10    else:
11        regex = r"(?:([1-9]+?.?)|([IVX]*.))?(\\s+)?(?: (Abstract)|(ABSTRACT))((\\s+)? *\\n?)(?P<
12        text>(?:.\\n)*?)^(?([1-9]+?.?)|([IVX]*.))?.s+(Introductio.n|INTRODUCTIO.N| A r t i c i a l ))|
13        Introductio.n\\n)"
14    matches = re.finditer(regex, text, re.MULTILINE)
15    # Parcours de tous les groupes pour debug
16    for matchNum, match in enumerate(matches, start=1):
17
18        #print ("Match {matchNum} was found at {start}-{end}: {match}".format(matchNum =
19        matchNum, start = match.start(), end = match.end(), match = match.group()))
20
21        for groupNum in range(0, len(match.groups())):
22            groupNum = groupNum + 1
23            #print ("Group {groupNum} found at {start}-{end}: {group}".format(groupNum =
24            groupNum, start = match.start(groupNum), end = match.end(groupNum), group = match.group(
25            groupNum)))
26
27    # Rechercher le texte correspondant a la regex
28    groups = re.compile(regex)
29    textIndex = groups.groupindex['text']
30    abstract_match = re.findall(regex, text, re.MULTILINE)
31    abstract = "N/A"
32    # Verifier si un resultat a ete trouve
33    if abstract_match:
34        # Afficher le texte extrait
35        abstract = match.group(textIndex)
36    return abstract.replace('\\n', ' ')

```

Listing 1: fonction getAbstract()

Une fois l'ensemble de nos expressions régulières écrites pour toutes les parties d'une publication, il ne restait plus qu'à écrire une méthode permettant de générer le fichier xml ou le fichier texte.

Cette étape était simple, il suffisait d'écrire dans un fichier les chaînes de caractères récupérées dans nos méthodes, au détail près qu'il fallait rajouter les balises XML dans le cas du fichier XML en sortie.

```

1 def writeXML(file_name, output_file_name, text, metadata, pdf):
2     """
3     Writes all the capital information in a .xml file with an XML layout
4     """
5     outputXML = "<article>\n"
6     outputXML+="<t<preamble> "+file_name+"</preamble>\n"
7     outputXML+="<t<titre>"+getTitle(metadata, text)+"</titre>\n"
8     outputXML+="<t<auteurs>\n"
9     auteursInfo = getAuthors(metadata, text, getTitle(metadata, text))
10
11    for i in auteursInfo:
12
13        outputXML+="<t\t<auteur>\n"
14        outputXML+="<t\t\t<name>"+i+"</name>\n"
15        outputXML+="<t\t\t\t<mail>"
16        outputXML+=auteursInfo[i]['mail']
17        outputXML+="</mail>\n"
18        outputXML+="<t\t\t\t<affiliation>"+auteursInfo[i]['affiliation']+"</affiliation>\n"
19        outputXML+="<t\t\t</auteur>\n"
20    outputXML+="<t</auteurs>\n"
21    if(file_name=="IPM1481.pdf"):
22        outputXML+="<t<abstract>"+getAbstract(pdf.pages[1].extract_text(), file_name)+"</
23        abstract>\n"

```

```

23 else:
24     outputXML+="\t<abstract>"+getAbstract(text)+"</abstract>\n"
25     outputXML+="\t<introduction>"+getIntroduction(text,file_name)+"</introduction>\n"
26     outputXML+="\t<discussion>"+getDiscussion(text)+"</discussion>\n"
27     outputXML+="\t<conclusion>"+getConclusion(text,file_name)+"</conclusion>\n"
28
29     outputXML+="\t<biblio>"+getBiblio(text,file_name)+"</biblio>\n"
30
31     outputXML+= "</article>"
32     if(output_file_name!=""):
33         fd = os.open(output_file_name,flags=os.O_RDWR | os.O_CREAT | os.O_TRUNC)
34         text = str.encode(outputXML)
35         lgtext = os.write(fd,text)
36         if(lgtext==0):
37             sys.stderr("Aucune donnees n'a pu etre extraite")
38         os.close(fd)
39     return outputXML

```

Listing 2: fonction writeXML()

### 3 Résultats

Plus le projet avançait plus les résultats que nous avons obtenus s'amélioraient. Nous vous

présenterons donc ici les résultats obtenus pour le dernier sprint du projet basés sur un corpus d'articles scientifiques de test. Voici les résultats obtenus sur ce dernier :

	Moyenne		acl2012	b0e5c43edf116c BLESS	C14-1212	Guy	InfoEmbeddings	IPM1481	L18	AI-Morality	SurveyTerm
Preamble	99,333	=	100	100	100	100	93,33	100	100	100	100
Titre	99,217	=	99,07	99,29	98,99	98,99	99,47	100	100	100	96,36
Auteurs	88,011	↗	99,4	93,37	93,4	96,27	99,56	92,95	94	90,43	66,05
Introduction	92,767	↗	99,85	98,85	99,78	93,9	99,87	96,72	59,99	99,93	100
Abstract	95,026	↗	99,68	99,66	99,67	99,7	99,86	82,59	99,2	99,64	99,01
Conclusion	88,298	↗	98,89	89,31	82,28	96,48	99,14	98,29	49,87	98,82	96,72
Bibliographie	89,766	↗	95,78	98,43	97,56	96,15	65,38	93,55	93,05	71,74	96,29
Total	93,081	↗	97,73	96,99	95,95	97,36	93,8	94,87	85,16	94,37	94,01
			=	↗	↗	↘	↗	↗	↗	=	=

Figure 2: Résultat finaux du logiciel

Nous pouvons observer que les cellules en vert sont les meilleurs résultats alors que ceux en blancs et en rouges sont moins bons. Ces scores décevants peuvent s'expliquer de plusieurs manières. Premièrement nos expressions régulières peuvent ne pas fonctionner en toutes circonstances et donc ne pas capturer complètement les parties de texte concernées. On a également retrouvé des différences notables entre les documents de

références permettant d'obtenir les scores et le contenu véritable des PDF qui devait être extrait, provoquant parfois de légères baisses au score final. Durant ce sprint, nos scores n'ont pas cessés d'augmenter, en grande partie en raison de l'amélioration de notre recherche des auteurs qui n'était jusqu'alors pas ou peu intégrée au code, et ce jusqu'à atteindre un score de 93%.

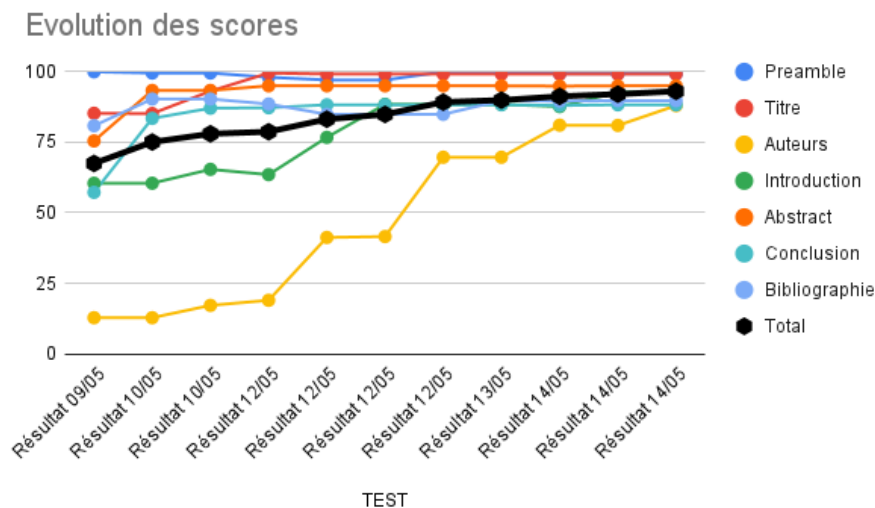


Figure 3: Evolution des scores

## 4 Conclusion

Notre système a su répondre aux attentes de notre client au vu des bons résultats que nous avons obtenus. Grâce à des ajustements durant le sprint final sur le corpus de test, nous avons optimisé le plus possible l'ensemble de nos expressions régulières afin de capturer un maximum d'information, et d'augmenter le taux de précision global au meilleur niveau possible. Ce cadre de la méthode agile SCRUM nous a permis lors de chaque sprint d'améliorer un peu plus notre méthode ainsi que nos scores sur l'ensemble des textes. Le projet nous a fait découvrir une nouvelle façon d'appréhender le travail de groupe. Cependant nous avons tout de même rencontré quelques difficultés tout au long du projet :

- Le changement entre le corpus d'entraînement

et celui de test qui comportait des PDF formatés de nouvelles manières

- Le contenu incorrect et très changeant des XML utilisé pour calculer le pourcentage qui nous ont fait changer très souvent nos méthodes à la fin du projet
- L'outil utilisé pour convertir nos PDF en texte qui n'était pas toujours très performant et introduisait des erreurs ou complications dans le texte à analyser

Le système produit à l'issue de ce projet est complet, efficace avec un score largement satisfaisant et ne contient que très peu de règles mises en place uniquement pour un fichier. On peut donc penser qu'il garderait une efficacité très suffisante en cas de généralisation et d'utilisation sur de nouveaux articles correctement formatés.