

PyMedExt, un couteau suisse pour le traitement des textes médicaux

William Digan, Antoine Neuraz, Alice Rogier,
David Baudoin, Nicolas Garcelon, Bastien Rance
Journée AFIA-TLH/ATALA - TAL & IA 2021

jeudi 4 février 2021

Prise en charge des patients



Dossier Patient Informatisé



Entrepôt de Données Cliniques



Entrepôt de données de l'Hôpital Européen
Georges Pompidou

Données intégrées depuis 2000
Données pour environ 1 million de patients

Concept	# observations
Items du DPI	176,524,243
Examens de biologie	155,173,140
Codes diagnostics	4,846,602
Prescriptions médicamenteuses	91,251,062

Grande variété de types de documents

Documents

- Comptes rendus de consultation
- Comptes rendus d'hospitalisation
- Comptes rendus opératoires
- Comptes rendus d'imagerie
- Comptes rendus d'anatomo-pathologie
- Lettres
- Ordonnances
- ...

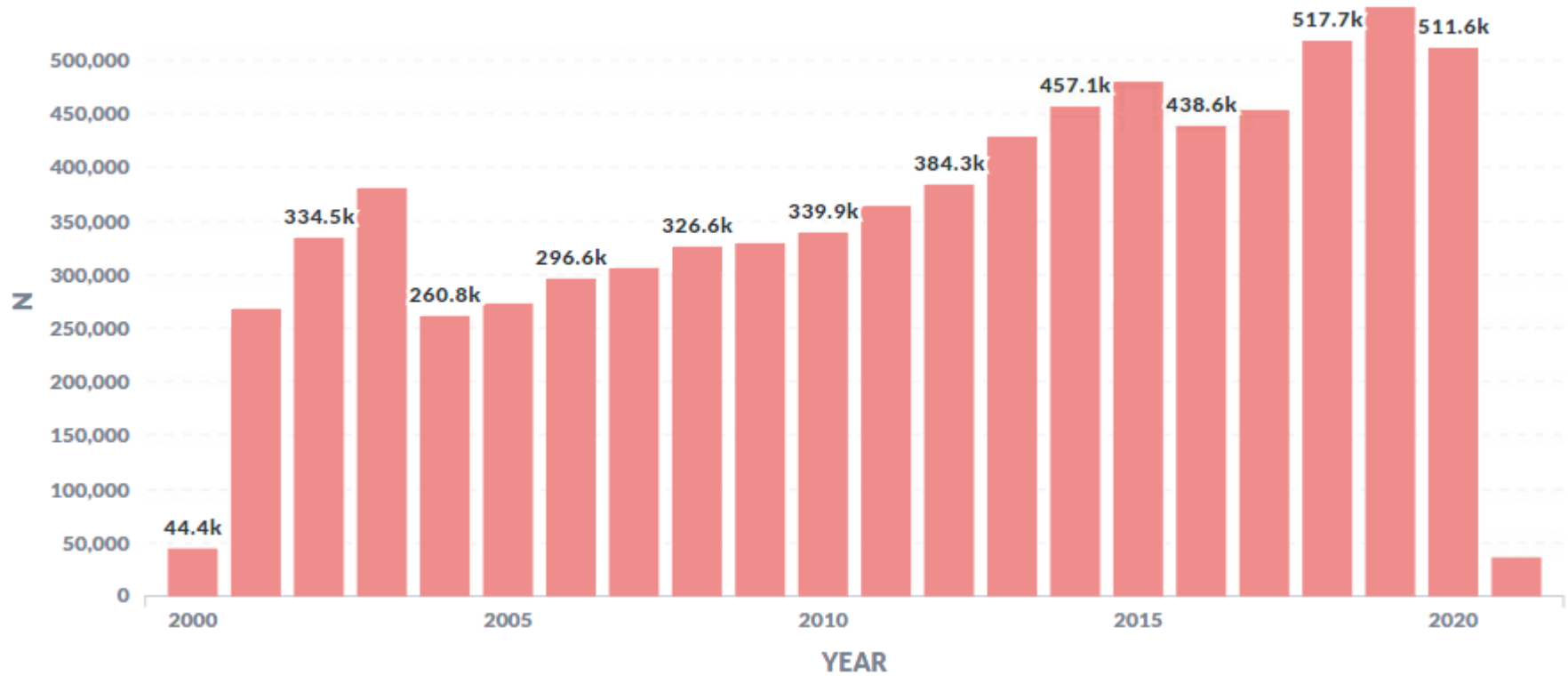
Et aussi

- Champs textes du dossier patient
- Transmissions infirmières
- ...

Concept	# observations
Item de questionnaires	176,524,243
Biologie	155,173,140
Transmissions infirmières	22,001,856
PMSI	4,846,602
Prescription médicamenteuses	91,251,062
Textes cliniques	5,343,224
Textes complémentaires	1,650,027

LES TEXTES DANS L'HÔPITAL

Statistiques de l'Hôpital Européen Georges Pompidou, AP-HP

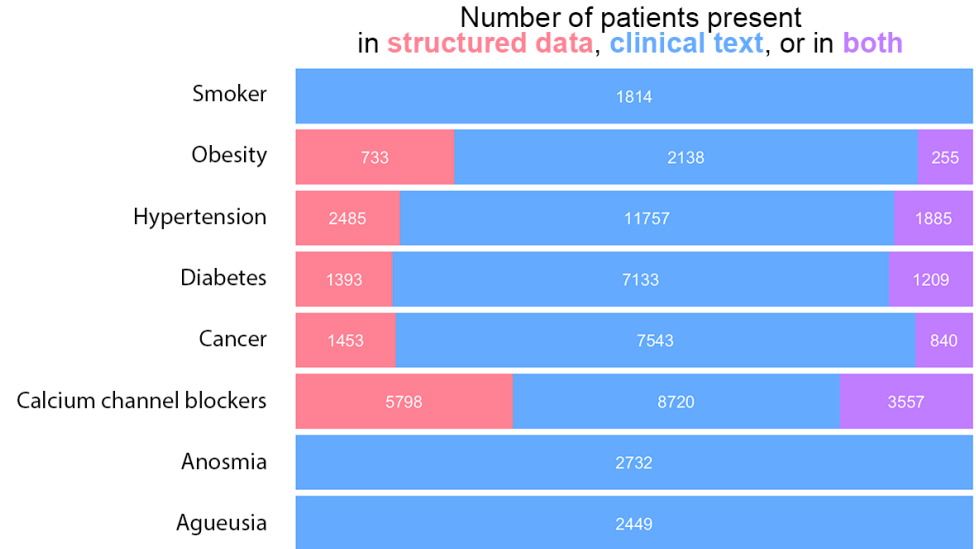


Plus de 1200 nouveaux documents chaque jour

L'IMPORTANCE DES TEXTES CLINIQUES

80%

De l'information médicale n'est présente que dans les textes



LE CIRCUIT DES DONNÉES

Prise en charge des patients



Dossier Patient Informatisé



Entrepôt de Données Cliniques



Traitement des textes



Recherche de patients éligibles pour des études



Extraction d'informations pour des recherches



Phénotypage à haut-débit



QUELLES RESSOURCES POUR TRAITER LES TEXTES



Nombreuses solutions intégrées (cTakes, CLAMP, GATE...)

Surtout en anglais

Souvent liées au formalisme d'annoteurs UIMA



Nombreux outils (Metamap, QuickUMLS, HeidelTime...)

Grande variété de langages de programmations

Grande variété de formats d'entrées et sorties

Des solutions de partage de *pipelines* et d'intégration

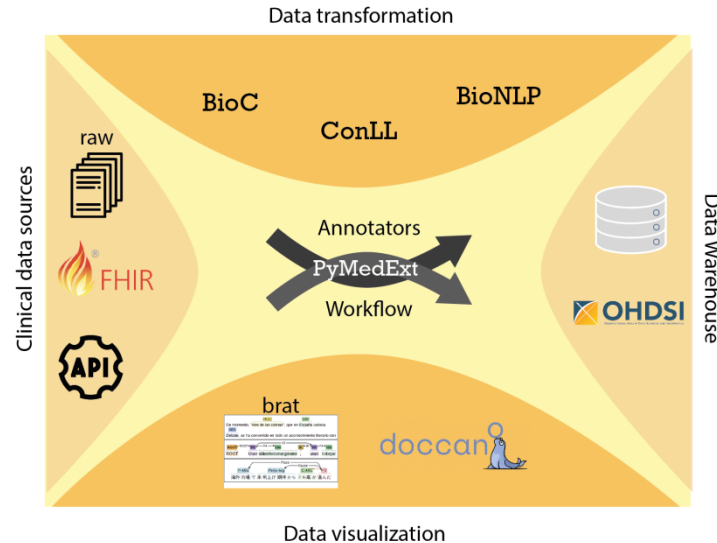


REPRODUCTIBILITÉ EN TRAITEMENT DES LANGUES

PYMEDEXT UNE PROPOSITION POUR FACILITER LA REPRODUCTIBILITÉ

Un couteau suisse pour le traitement des langues cliniques

- Un **format de représentation** des documents et des annotations (extension de BioC)
- Des **convertisseurs** de formats
- Des **exports** vers des formats standards ou outils d'annotation
- Un format de représentation **d'annoteurs**
- Un moteur de **pipelines simples**



Travail démarré pendant la crise Covid. Collaboration de toute l'équipe (PI Antoine Neuraz)

- Utilisé pour l'annotation des comptes rendus cliniques à l'AP-HP pendant la première vague
- Outils d'échange dans l'équipe

https://github.com/equipe22/pymedext_core

Le format PyMedExt

Document

Annotation

ID

Type

Valeur

Source

SourceID

Span(début,fin)

Le format PyMedExt



Le format PyMedExt



Philosophie de PyMedExt pour les annotateurs :

Structure **simple**

- Héritage d'une classe python
- Deux fonctions à compléter

Deux possibilités pour la **création de nouveaux annotateurs** PyMedExt :

- Développement en python natif
- Encapsulation de codes existants

AJOUTER UN ANNOTATEUR EN UIMA

1 Type System holds CAS Feature Structure representation(s)

2 The Analysis Engine Descriptor

3 Generate Java classes

4 The Annotator Class

```

22 <typesystemDescription xmlns="http://uima.apache.org/resourceSpecifier">
23   <name>SetTypeSystem</name>
24   <description>Type System Definition for SETH</description>
25   <vendor>Apache Software Foundation</vendor>
26   <types>
27     <typeDescription>
28       <name>de.dfxi.it.finds.finds.MutationAnnotation</name>
29       <description>Mutation annotation created with SETH</description>
30       <superTypeNames>uima.cas.Annotation</superTypeNames>
31       <features>
32         <featureDescription>
33           <name>Type</name>
34           <description>Metadata type</description>
35           <rangeTypeNames>uima.cas.String</rangeTypeNames>
36         </featureDescription>
37         <featureDescription>
38           <name>Features</name>
39           <description>Mutation</description>
40           <rangeTypeNames>uima.cas.String</rangeTypeNames>
41         </featureDescription>
42         <featureDescription>
43           <name>MetaDesc</name>
44           <description>MetaDesc</description>
45           <rangeTypeNames>uima.cas.String</rangeTypeNames>
46         </featureDescription>
47         <featureDescription>
48           <name>Position</name>
49           <description>Metadata position</description>
50           <rangeTypeNames>uima.cas.String</rangeTypeNames>
51         </featureDescription>
52         <featureDescription>
53           <name>Tool</name>
54           <description>Annotation tool</description>
55           <rangeTypeNames>uima.cas.String</rangeTypeNames>
56         </featureDescription>
57         <featureDescription>
58           <name>Types</name>
59           <description>Types</description>
60           <rangeTypeNames>uima.cas.String</rangeTypeNames>
61         </featureDescription>
62       </features>
63     </typeDescription>
64   </types>
65 </typesystemDescription>

```

```

22 <!-- Descriptor for the MutationAnnotator - ->
23 <analysisEngineDescription xmlns="http://uima.apache.org/resourceSpecifier">
24   <frameworkImplementation>org.apache.uima.java</frameworkImplementation>
25   <primitive>true</primitive>
26   <annotatorImplementationName>
27     de.dfxi.it.finds.finds.MutationAnnotator
28   </annotatorImplementationName>
29   <analysisEngineMetadata>
30     <name>SETH Mutation Annotator</name>
31     <description>A mutation annotator using SETH (SNP Extraction Tool for Human Variations).</description>
32     <version>1.0</version>
33     <creator>Deutsches Forschungszentrum für Künstliche Intelligenz (DFKI)</creator>
34   </analysisEngineMetadata>
35   <typeSystemDescription>
36     <imports>
37       <import location="SetTypeSystem.xml"/>
38     </imports>
39   </typeSystemDescription>
40   <capabilities>
41     <capability>
42       <input>/>
43       <type allAnnotatorFeatures="true">de.dfxi.it.finds.finds.MutationAnnotation</type>
44     </output>
45     </capability>
46   </capabilities>
47   <analysisEngineMetadata>
48     <analysisEngineMetadata>

```

```

<dependency>
<groupId>org.apache.uima</groupId>
<artifactId>jcasgen-maven-plugin</artifactId>
<version>3.0.0</version>
</dependency>

```

```

1 // Import classes from org.apache.uima
2 import org.apache.uima.analysisengine.impl.*;
3 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
4 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
5 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
6 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
7 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
8 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
9 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
10 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
11 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
12 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
13 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
14 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
15 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
16 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
17 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
18 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
19 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
20 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
21 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
22 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
23 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
24 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
25 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
26 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
27 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
28 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
29 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
30 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
31 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
32 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
33 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
34 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
35 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
36 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
37 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
38 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
39 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
40 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
41 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
42 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
43 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
44 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
45 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
46 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
47 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
48 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
49 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
50 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
51 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
52 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
53 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
54 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
55 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
56 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
57 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
58 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
59 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
60 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
61 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
62 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
63 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
64 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
65 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
66 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
67 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
68 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
69 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
70 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
71 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
72 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
73 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
74 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
75 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
76 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
77 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
78 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
79 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
80 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
81 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
82 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
83 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
84 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
85 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
86 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
87 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
88 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
89 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
90 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
91 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
92 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
93 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
94 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
95 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
96 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
97 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
98 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
99 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;
100 import org.apache.uima.analysisengine.impl.analysisenginefactory.*;

```

5 Calling the Annotator

```

47 public class FIMDA {
48
49   private AnalysisEngine ae;
50
51
52   FIMDA() throws IOException, InvalidXMLException, ResourceInitializationException {
53     //get Resource Specifier from XML File
54     XMLInputSource in;
55     URL url = getClass().getResource("/resources/desc/MutationAnnotator.xml");
56     // important: use URL to create the XMLInputSource, otherwise it wont work from the jar package
57     in = new XMLInputSource(url);
58     ResourceSpecifier specifier = UIMAFramework.getXMLParser().parseResourceSpecifier(in);
59
60
61     //creates AE here
62     ae = UIMAFramework.produceAnalysisEngine(specifier);
63   }
64 }

```

Définition de l'annotateur

(héritage de la classe annotateur)

Implémentation de la fonction *annotate_function*

```
class findMatches(annotators.Annotator):
    """
    Annotator based on linux grep to search regex from a source file
    """
    def __init__(self, key_input, key_output, ID, findValues ):
        """FIXME! initialize the annotator

        :param key_input: input ['raw_text']
        :param key_output: Annotation type here "Liposarcom.V0"
        :param ID: regex_fast.version
        :param findValues: "list of value to identify in the text"
        :returns:
        :rtype:

        """
        super().__init__(key_input, key_output, ID)
        self.findValues=findValues
```

```
def annotate_function(self, _input):
    """ main annotation function
    :param _input: in this case raw_text
    :returns: a list of annotations
    :rtype:
    """
    logger.debug(_input)
    inp = self.get_key_input(_input,0)[0]
    annotationsList=[]
    for thisValue in self.findValues:
        #result = [i.start() for i in re.finditer(thisValue, inp.value.lower())]
        for i in re.finditer(thisValue, inp.value.lower()):
            matchPos=i.start()
            if matchPos is not []:
                logger.debug("ok go in loop")
                logger.debug(matchPos)
                ID = str(uuid.uuid1())
                annotationsList.append(annotators.Annotation(type= self.key_output,
                    value=thisValue, #thisMatch,
                    span=(int(matchPos), int(matchPos)+len(thisValue)),
                    source=self.ID,
                    isEntity=True,
                    ID=ID,
                    source_ID = inp.ID))

    logger.debug(annotationsList)
    return(annotationsList)
```

Définition de l'annotateur

(héritage de la classe annotateur)

Implémentation de la fonction *annotate_function*

```
class findMatches(annotators.Annotator):
    """
    Annotator based on linux grep to search regex from a source file
    """
    def __init__(self, key_input, key_output, ID, findValues ):
        """FIXME! initialize the annotator

        :param key_input: input ['raw_text']
        :param key_output: Annotation type here "Liposarcom.V0"
        :param ID: regex_fast.version
        :param findValues: "list of value to identify in the text"
        :returns:
        :rtype:

        """
        super().__init__(key_input, key_output, ID)
        self.findValues=findValues
```

```
def annotate_function(self, _input):
    """ main annotation function
    :param _input: in this case raw_text
    :returns: a list of annotations
    :rtype:
    """
    logger.debug(_input)
    inp = self.get_key_input(_input,0)[0]
    annotationsList=[]
    for thisValue in self.findValues:
        #result = [i.start() for i in re.finditer(thisValue, inp.value.lower())]
        for i in re.finditer(thisValue, inp.value.lower()):
            matchPos=i.start()
            if matchPos is not []:
                logger.debug("ok go in loop")
                logger.debug(matchPos)
                ID = str(uuid.uuid1())
                annotationsList.append(annotators.Annotation(type= self.key_output,
                    value=thisValue, #thisMatch,
                    span=(int(matchPos), int(matchPos)+len(thisValue)),
                    source=self.ID,
                    isEntity=True,
                    ID=ID,
                    source_ID = inp.ID))

            logger.debug(annotationsList)
    return(annotationsList)
```

ANNOTATEURS PYMEDEXT

1

```
class findMatches(annotators.Annotator):
    """
    Annotator based on linux grep to search regex from a source file
    """
    def __init__(self, key_input, key_output, ID, findValues ):
        """FIXME! initialize the annotator

        :param key_input: input ['raw_text']
        :param key_output: Annotation type here "Liposarcom.V0"
        :param ID: regex_fast.version
        :param findValues: "list of value to identify in the text"
        :returns:
        :rtype:
        """
        super().__init__(key_input, key_output, ID)
        self.findValues=findValues
```

```
def annotate_function(self, _input):
    """ main annotation function
    :param _input: in this case raw_text
    :returns: a list of annotations
    :rtype:
    """
    logger.debug(_input)
    inp = self.get_key_input(_input,0)[0]
    annotationsList=[]
    for thisValue in self.findValues:
        #result = [i.start() for i in re.finditer(thisValue, inp.value.lower())]
        for i in re.finditer(thisValue, inp.value.lower()):
            matchPos=i.start()
            if matchPos is not []:
                logger.debug("ok go in loop")
                logger.debug(matchPos)
                ID = str(uuid.uuid1())
                annotationsList.append(annotators.Annotation(type= self.key_output,
                                                            value=thisValue, #thisMatch,
                                                            span=(int(matchPos), int(matchPos)+len(thisValue)),
                                                            source=self.ID,
                                                            isEntity=True,
                                                            ID=ID,
                                                            source_ID = inp.ID))
            logger.debug(annotationsList)
    return(annotationsList)
```

2

Appel de l'annotateur

```
demoAnnotator = findMatches(key_input = ['raw_text'],
                             key_output = 'Liposarcom.V0',
                             ID = "demoreiter", findValues = ["liposarcome"])

# add all your annotators in a list
annotatorsList =[demoAnnotator]
# annotate your document
LetterPyMedExt.annotate(annotatorsList)
```

Pipelines d'annotation simples

La gestion de pipelines complexes sera faite en utilisant des outils dédiés hors de PyMedExt

EXEMPLE D'ANNOTATEURS AU FORMAT PYMEDEXT

Développement en python natif

- Regex
- Médicaments
- Necker/Imagine contexte

Encapsulation de codes existants

- QuickUMLS
- HeidelTime
- Romedi

Prochainement rendus disponibles
https://github.com/equipe22/pymedext_core

CONVERSION DE FORMATS

Conversions

- Des **convertisseurs** de formats
- Des **exports** vers des formats standards ou outils d'annotation

Entrées

PyMedExt

Text

FHIR

BioC

CONVERSION DE FORMATS

Conversions

- Des **convertisseurs** de formats
- Des **exports** vers des formats standards ou outils d'annotation

Entrées

Sorties

PyMedExt

PyMedExt-json

Text

FHIR

BioC

Standards
informatique
médicale

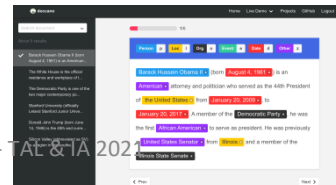
Annotation
Visualisation

OMOP



BRAT

Doccano



CONVERSION DE FORMATS

Conversions

- Des **convertisseurs** de formats
- Des **exports** vers des formats standards ou outils d'annotation

Entrées

Sorties

En cours de développement

PyMedExt

PyMedExt-json

BioC

Text

CoNLL

FHIR

OMOP



FHIR

BioC

Standards
informatique
médicale

Annotation
Visualisation

BRAT

Doccano



EXPORT VERS BRAT

1

```
class findMatches(annotators.Annotator):
    """
    Annotator based on linux grep to search regex from a source file
    """
    def __init__(self, key_input, key_output, ID, findValues ):
        """FIXME! initialize the annotator

        :param key_input: input ['raw_text']
        :param key_output: Annotation type here "Liposarcom.V0"
        :param ID: regex_fast.version
        :param findValues: "list of value to identify in the text"
        :returns:
        :rtype:
        """
        super().__init__(key_input, key_output, ID)
        self.findValues=findValues
```

```
def annotate_function(self, _input):
    """ main annotation function
    :param _input: in this case raw_text
    :returns: a list of annotations
    :rtype:
    """
    logger.debug(_input)
    inp = self.get_key_input(_input,0)[0]
    annotationsList=[]
    for thisValue in self.findValues:
        #result = [i.start() for i in re.finditer(thisValue, inp.value.lower())]
        for i in re.finditer(thisValue, inp.value.lower()):
            matchPos=i.start()
            if matchPos is not []:
                logger.debug("ok go in loop")
                logger.debug(matchPos)
                ID = str(uuid.uuid1())
                annotationsList.append(annotators.Annotation(type= self.key_output,
                                                            value=thisValue, #thisMatch,
                                                            span=(int(matchPos), int(matchPos)+len(thisValue)),
                                                            source=self.ID,
                                                            isEntity=True,
                                                            ID=ID,
                                                            source_ID = inp.ID))
            logger.debug(annotationsList)
    return(annotationsList)
```

2

Appel de l'annotateur

```
demoAnnotator = findMatches(key_input = ['raw_text'],
                             key_output = 'Liposarcom.V0',
                             ID = "demoreiter", findValues = ["liposarcome"])

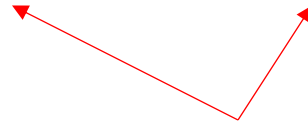
# add all your annotators in a list
annotatorsList =[demoAnnotator]
# annotate your document
LetterPyMedExt.annotate(annotatorsList)
```

```
pymedext.brat.savetobrat(LetterPyMedExt,path)
```

CONVERSION DE FORMATS – EN LIGNE DE COMMANDE

```
pymedext -i demo.txt --itype txt -otype pymedext
```

```
pymedext -i MEDLINE_train_bioc --itype biocjson -otype pymedext
```



Format explicite

spaCy

spaCy

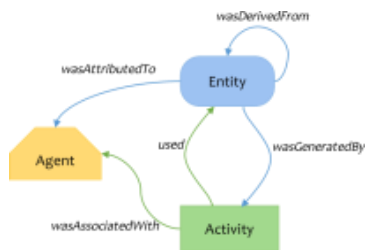
Large communauté d'utilisateurs

Optimisation code Python et Cython

Tokenisation (Hypothèse forte)

Travail en cours pour :

- intégrer des annotateurs spaCy en PyMedExt
- Rendre compatible les modèles de documents



Reproductibilité

- PROV-O ontologie pour l'expression de la provenance
- Travail sur une containerisation des outils externes et de PyMedExt

CONCLUSION

Projet récent

- **Utilisation** dans les hôpitaux Necker/Institut Imagine et HEGP
- Simplifie le partage de ressources
- Prise en main rapide

Encore du travail nécessaire, notamment pour éviter la redondance avec d'autres projets

Ressources distribuées sur le github du projet

Les testeuses et testeurs, contributrices et contributeurs *sont les bienvenues* !

PyMedExt en action :

- **Prochaine présentation** avec Antoine Neuraz !

From the Noun Project

- data warehouse by remmachenasreddine
- Mining by Firza Alamsyah
- text file by Wing
- Electronic Health Record by Bold Yellow
- extraction by Eucalyp
- Hospital by ProSymbols