

Documentation LaToE 2.0 :

Ce document vous donnera une idée du procédé d'identification de structures mis en place par LAPDFTEXT, et la manière dont celui-ci a été altéré.

La classe centrales de la librairie sont les classes **RuleBasedParser**, qui gère tout le parsing, **Document** qui représente le PDF parsé et **RuleBasedChunkClassifier** qui permet la labellisation des blocs. Voici un aperçu des classes principales utilisées. Il y en a d'autres dans la librairie mais on n'a pas interagi de manière directe avec elles.

LayoutAwarePDFText	
Source Packages	
<default package>	
Interface.java	Classes ajoutées à LAPDFTEXT (interface graphique, générateur de règles, etc)
Rule.java	
RuleGenerator.java	
edu.isi.bmkeg.pdf	
DocumentInformation.java	Classes non utilisées, n'interviennent pas directement dans le parsing
DocumentInformation_Type.java	
edu.isi.bmkeg.pdf.classification	
Classifier.java	Interface générique pour la classification
edu.isi.bmkeg.pdf.classification.ruleBased	
RuleBasedChunkClassifier.java	Classe permettant la labellisation, implémente l'interface ci-dessus
edu.isi.bmkeg.pdf.extraction	
Extractor.java	Classes liées à l'extraction des wordblocks
JPedaExtractor.java	
JPedaPageImageExtractor.java	
edu.isi.bmkeg.pdf.extraction.exceptions	
AccessException.java	Exceptions spécifiques à la librairie
EmptyPDFException.java	
EncryptionException.java	
InvalidPopularSpaceValueException.java	
edu.isi.bmkeg.pdf.features	
ChunkFeatures.java	Classes permettant de définir des propriétés pour les blocs, propriétés utiles pour la classification
HorizontalSplitFeature.java	
WordFeatures.java	
edu.isi.bmkeg.pdf.model	
Block.java	Interfaces génériques pour les types de blocs. Document en revanche est une classe à part.
ChunkBlock.java	
Document.java	
PageBlock.java	
WordBlock.java	
edu.isi.bmkeg.pdf.model.RTree	
RTChunkBlock.java	Implémentation de ces interface par rapport à une logique RTree
RTDummyProcedure.java	
RTModelFactory.java	
RTPageBlock.java	
RTProcedure.java	
RTSpatialEntity.java	
RTSpatialRepresentation.java	
RTWordBlock.java	
edu.isi.bmkeg.pdf.model.factory	
AbstractModelFactory.java	Contient les définitions des méthodes de création de blocs
edu.isi.bmkeg.pdf.model.ordering	
SpatialOrdering.java	Méthodes de classement spatial des blocs (horizontal. vertical. etc)
edu.isi.bmkeg.pdf.model.spatial	
SpatialEntity.java	Classes permettant la représentation spatiale d'un bloc. Tous les blocs du package Model étendent l'une de ces deux classes.
SpatialRepresentation.java	
edu.isi.bmkeg.pdf.parser	
Parser.java	Interface et implémentation du parseur. Classe qui permet le lancement du parsing.
RuleBasedParser.java	

Déroulement du parsing en bref :

Afin de parser un PDF, on instancie la classe **RuleBasedParser**. Pour cela, on utilise le constructeur en donnant en paramètre une nouvelle instance de **ModelFactory**. Cette classe sert juste à définir la manière avec laquelle vont être créés les différents blocs.

```
obj = new RuleBasedParser(new RTModelFactory());
```

Une fois le parseur opérationnel, on peut lancer le parsing en utilisant la fonction **parse()** de la classe **RuleBasedParser**. Cette fonction prend en paramètre le chemin du fichier, et renvoie une instance de Document qui correspond à notre PDF parsé.

Cette fonction **parse()** fait elle-même appel à une autre classe, la classe **JPedalExtractor**. Cette classe possède plusieurs fonctions essentielles : **decodeFile()** qui permet d'extraire une liste de tous les mots du document, et **hasNext()** qui renvoie les mots un par un. Il est important de noter que cette classe fonctionne page par page.

Cette dernière classe est notre frontière dans ce projet car pour effectuer ses opérations elle fait appel à de nombreuses classes de librairies externes non modifiables, comme **PDFDecoder** et **PdfGroupingAlgorithms**. Ces classes proviennent pour la majorité de la librairie **JPedal**.

Revenons à la fonction **parse()**. Celle-ci va donc, grâce à la fonction **hasNext()** de la classe **JPedalExtractor**, recevoir des entités **WordBlock** correspondant à tous les mots du PDF. Ces entités contiennent diverses informations comme les coordonnées, la largeur, la hauteur, la police, la taille de police, le texte...

Le but du parsing est de former à partir de ces **WordBlocks** des **ChunkBlocks**. Ces chunks vont ensuite correspondre à des paragraphes, des titres, etc. Donc pour chaque page, la fonction **parse()** va extraire tous les mots présents dans celle-ci, et à partir d'un mot, suivre l'algorithme suivant :

- Si le mot est le premier détecté, on crée un chunkblock qui ne contient que celui-ci
- A chaque nouveau mot, on vérifie si le mot est en contact avec le chunkblock précédent. Si c'est le cas, on fusionne le chunk et le mot. Sinon, on crée un nouveau chunk pour le nouveau mot.
- La notion de « contact » est définie par deux variables permettant une tolérance : NorthSouth et EastWest. Ces variables définies grâce à des statistiques sur le document déterminent une valeur moyenne d'écart horizontal et vertical entre les mots. Nous avons d'ailleurs modifié cette partie-là, en introduisant des multiplicateurs pour les variables. Ainsi, nous pouvons augmenter ou réduire la précision du parsing à volonté.

A la fin du parsing, une instance de Document est créée, qui renvoie sur une liste de pages (PageBlock), qui possèdent chacune une liste de chunkBlocks correspondant aux blocs de la page.

Déroulement de la labellisation :

Au cœur du processus de labellisation on trouve trois éléments : la classe **RuleBasedChunkClassifier**, qui permet de lancer l'analyse, le fichier de règles **rules.drl**, et la classe **ChunkFeatures** qui permet de définir des propriétés utilisables dans le fichier de règles.

Le fichier de règles est structuré suivant une syntaxe drools. Les règles ont la forme qui suit :

<pre>rule "item" no-loop activation-group "blockClassification" salience 3 when ChunkFeatures(item==true) then chunk.setType(chunk.TYPE_ITEM); end</pre>	<div>Propriétés de la règle salience => priorité</div> <div>Conditions</div> <div>Conséquences</div>
--	---

Ici, on se sert d'une fonction définie dans ChunkFeatures comme condition :

```
public boolean isItem() {
    String aux = chunk.getchunkText();
    aux=aux.replaceAll("<[^>]*>", "");
    aux=aux.replaceAll("\\s", "");
    if (aux.matches("[0-9]+|([a-z]|[A-Z]))[-.:].*")) {
        return true;
    } else {
        return false;
    }
}
```

Il est possible de définir n'importe quel type de propriétés de la même manière. Il suffit de créer une fonction qui a comme nom « get + nom de la propriété » ou « is + nom de la propriété » si la fonction crée ne renvoie que des booléens. Par exemple, la fonction `isItem()` ci-dessus détecte si le bloc en question est un item de liste.

On remplit donc le fichier drools suivant les blocs que l'on souhaite détecter, et on crée en parallèle des propriétés dans **ChunkFeatures** si besoin. Une fois terminé, on instancie la classe **RuleBasedChunkClassifier** comme suit :

```
RuleBasedChunkClassifier classifieur = new RuleBasedChunkClassifier(rules, new RTModelFactory());
```

Le constructeur prend en paramètre le chemin du fichier de règles (par défaut, le fichier se trouve dans `\src\main\resources\rules`) ainsi qu'une instance de **RTModelFactory**, que l'on crée. Une fois la classe instanciée, on lance la labellisation grâce à la fonction `classify()`. Cette fonction prend en paramètre la liste de chunks que l'on souhaite classer.

Le label des blocs est stocké dans la variable **type** des chunks. Sa valeur par défaut est « unclassified ».

Interface graphique :

The screenshot shows the main interface of the PDF parsing tool. At the top, there's a file path input field and a 'Parcourir' (Browse) button. Below this are sliders for 'Tolérance X' and 'Tolérance Y', both set to 50%. There are buttons for 'Préc.' (Previous), 'Charger le PDF' (Load PDF), and 'Suiv.' (Next). Checkboxes for 'Afficher mots' (Show words) and 'Sensib. Max.' (Max. sensitivity) are also present.

The central area displays a visual representation of the PDF content with various blocks labeled. Callout boxes point to specific features:

- Tolérances PENDANT le parsing**: Points to the tolerance sliders.
- Chemin du pdf**: Points to the file path input field.
- Rend le parsing le plus précis possible**: Points to the 'Sensib. Max.' checkbox.
- Analyse APRES le parsing : permet de fusionner certains chunks pour améliorer la précision. Possibilité de fusionner les blocs horizontalement ou verticalement, suivant une certaine tolérance.**: Points to the 'Analyse' tab and the 'Fusion verticale' and 'Fusion horizontale' checkboxes.
- Page précédente**: Points to the 'Préc.' button.
- Analyse le pdf suivant les tolérances X et Y et l'affiche**: Points to the 'Charger le PDF' button.
- Page suivante**: Points to the 'Suiv.' button.
- Affiche les wordblock**: Points to the 'Afficher mots' checkbox.
- Activer ou non la labellisation, spécifier le chemin du fichier de règles**: Points to the 'Labellisation' section and the 'rules.drl' file path.
- Détecte les blocs en dehors des marges, permet de les supprimer**: Points to the 'Détection des blocs en dehors des marges' checkbox.
- Analyse les blocs qui se poursuivent sur plusieurs pages**: Points to the 'Analyse des blocs qui se poursuivent sur plusieurs pages' checkbox.
- Analyse les relations de coordination/subordination et les affiche**: Points to the 'Analyse des relations de coordination/subordination et les affiche' checkbox.
- Effectue tous les traitements sur le PDF fourni au départ**: Points to the 'Traitement complet' button.
- Génération de plan suivant les relations hiérarchiques**: Points to the 'Générer plan' button.
- Sauvegarde/Chargement de paramètres et d'annotation**: Points to the 'Sauvegarder l'annotation' and 'Charger une annotation' buttons.
- Calcul de précision manuel VS automatique**: Points to the 'Calcul de précision' button.

On the right side, there's a sidebar with tabs for 'Analyse' and 'Clustering'. The 'Analyse' tab is active, showing various settings and buttons:

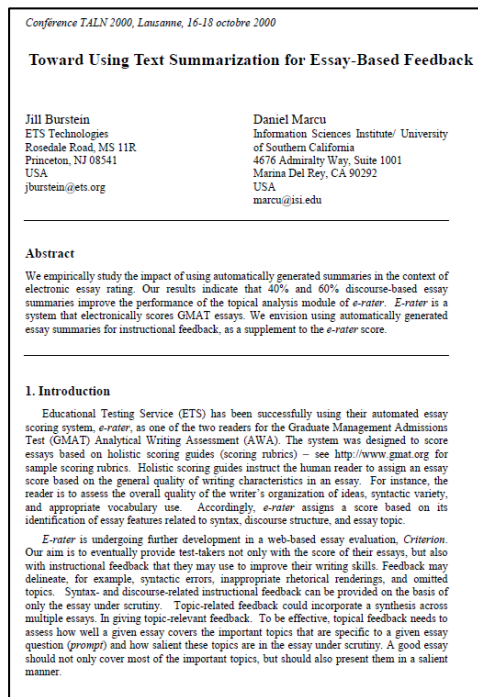
- Paramètres de l'analyse**: Includes checkboxes for 'Analyse post-parsing', 'Fusion verticale', 'Fusion horizontale', and 'Fusion fin de paragraphe'. There are input fields for '0.2' and '0.2'.
- Labellisation**: Includes a checkbox for 'Labelliser' and a text field for the rules file path ('/src/main/resources/rules/rules.drl').
- Redéfinition des bords**: Includes buttons for 'Analyser', 'Rogner', and 'RAZ'.
- Traitement complet**: A large red button.
- Générer plan**: A button.
- Annoter manuellement**: A button.
- Sauvegarder l'annotation**: A button.
- Charger une annotation**: A button.
- Calcul de précision**: A button.
- 0%**: A progress indicator.
- Log**: A text area showing the progress of the parsing process.

Le bouton « Traitement complet » effectue toutes les actions importantes sur le PDF fourni au départ.

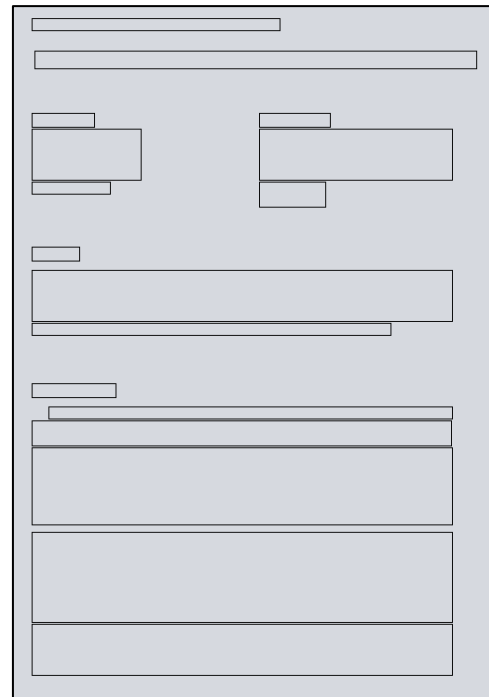
Options utilisables avant et pendant le parsing :

On peut sélectionner le PDF grâce au bouton « parcourir » ou entrer à la main un chemin valide. On peut ensuite mieux définir la précision lors du parsing : ces options modifient la marge avec laquelle un mot va être ajouté à un chunk (cf explication du parsing). Les précisions vont de 0 à 100%, 0

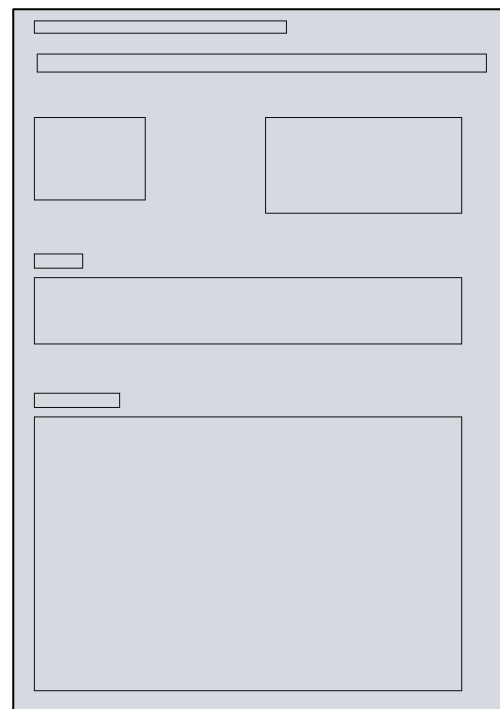
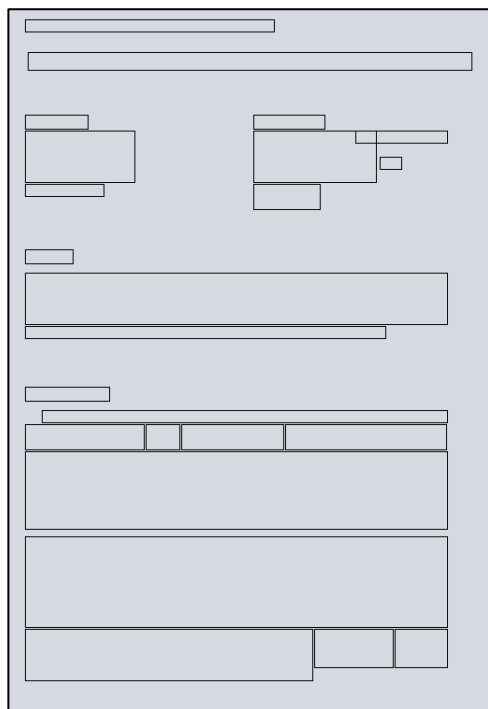
signifiant qu'un mot n'acceptera que ses voisins direct et 100 correspondant à la précision de base de la librairie LAPDFTEXT. Les tolérances sont de base à 50%. Une fois ceci fait, on peut lancer le parsing en cliquant sur « Charger le PDF ». Cette opération peut prendre plusieurs secondes. Une fois le PDF



chargé, les blocs s'affichent dans l'interface, comme suit :



Voici ce qui se passe avec des tolérances à 0% et à 100% :

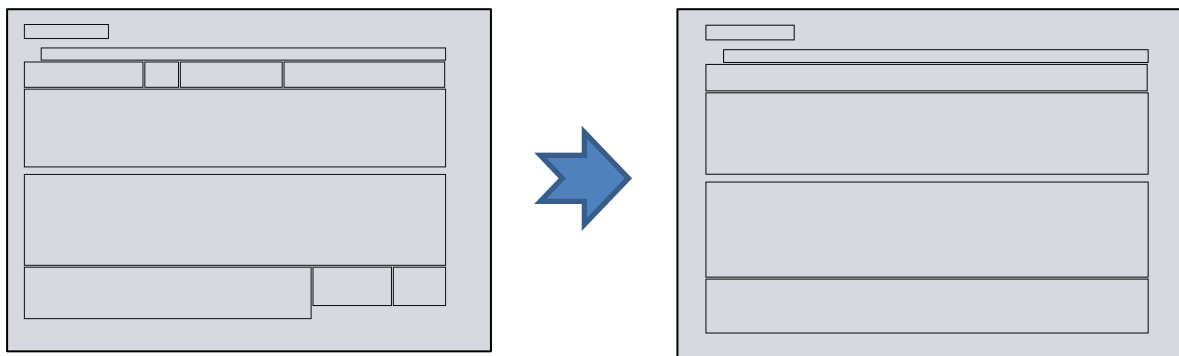


On constate donc bien l'impact des tolérances. Il est à noter que l'on peut définir des tolérances horizontales et verticales différentes.

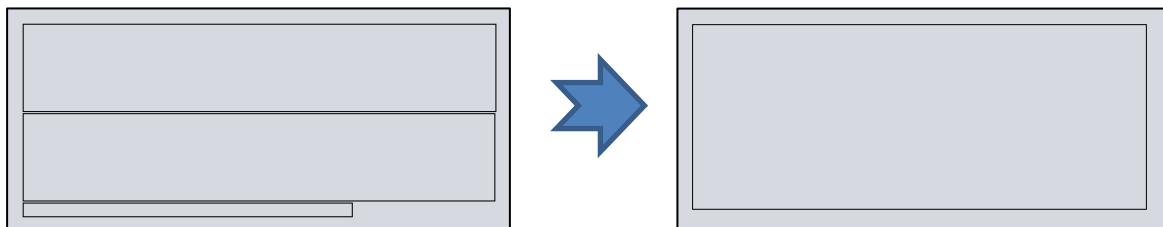
Analyse post-parsing :

L'analyse post-parsing a été intégralement ajoutée dans la librairie, après avoir constaté que le parsing était erroné quelques fois. Celle-ci agit après le parsing, en parcourant tous les blocs générés et en fusionnant ceux que l'on considère « trop proches » suivant certains critères. L'utilisateur peut choisir quelles fusions mettre en place (horizontale ou verticale) et la tolérance de ces fusions. Il n'y a pas de limite supérieure de tolérance, l'utilisateur est donc libre de gérer la fusion.

Exemple : fusion horizontale



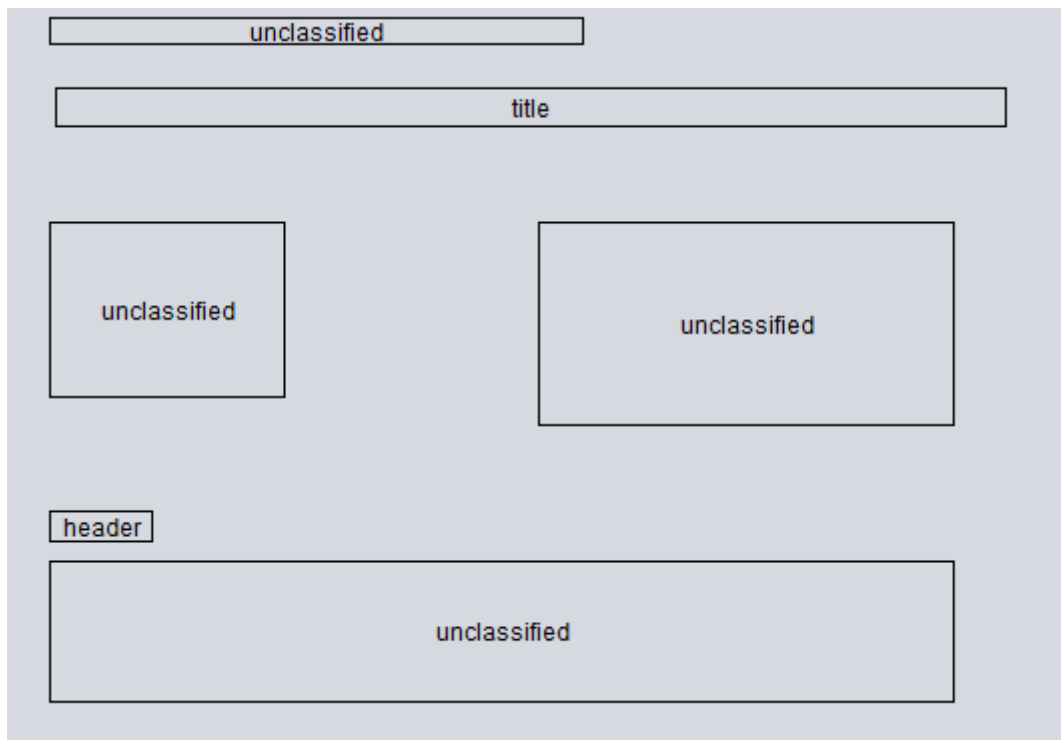
Il y a une troisième option dans l'analyse : la fusion paragraphe. En effet, il arrive que les alinéas en début ou en fin de paragraphe empêchent le parsing correct de celui-ci. L'analyse permet donc de régler ce problème.



Note : la fusion de paragraphe est bien distincte de la fusion verticale. Cette dernière n'effectue aucune vérification, et fusionne tous les blocs répondant aux critères de manière verticale. La fusion paragraphe elle n'agit que sur les blocs qu'elle détecte comme étant des débuts ou fins de paragraphes.

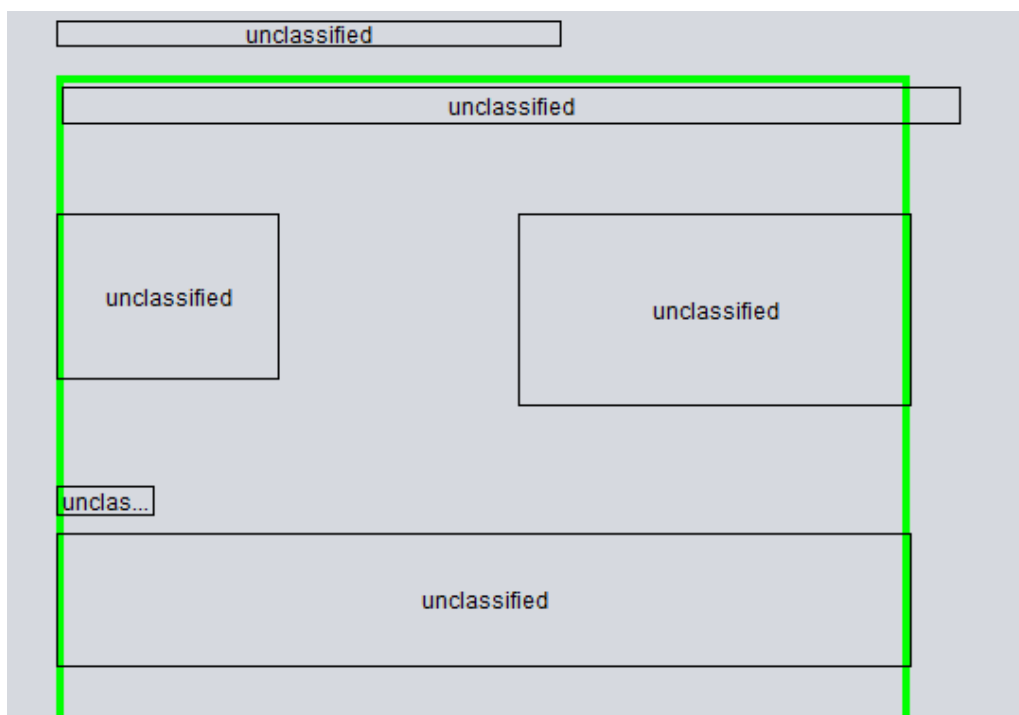
Labellisation :

L'utilisateur peut spécifier un fichier de règles différent, et activer ou non la labellisation grâce à la checkbox. Les types ainsi définis sont affichés au centre des blocs.



Redéfinition des bords :

L'utilisateur a trois options ici. L'option « **Analyser** » qui permet d'effectuer une étude statistique sur le document qui définit les marges de celui-ci. Ces marges sont affichées en vert.

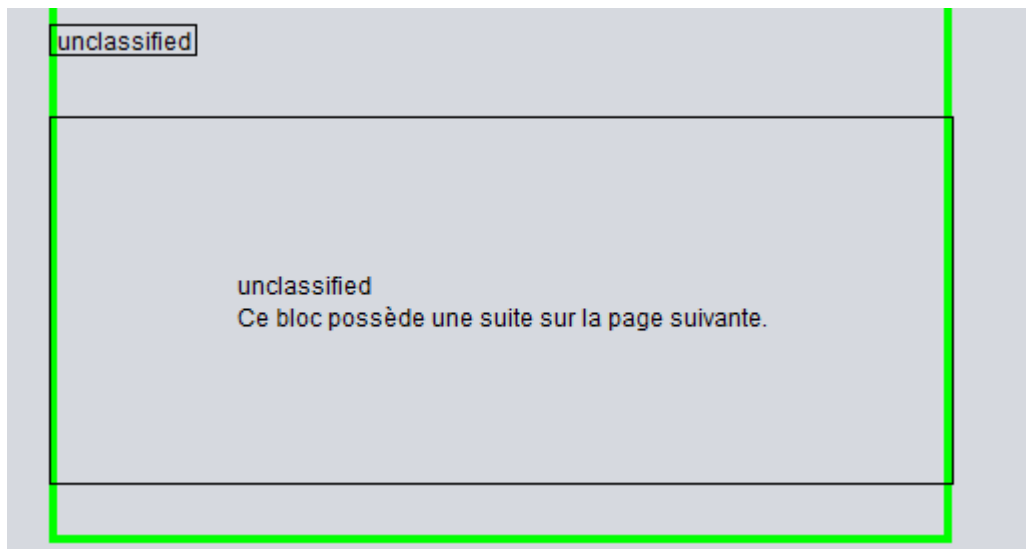


Ensuite, l'utilisateur peut sélectionner « **Rogner** » qui va supprimer les blocs en dehors des marges. Cette option est utile pour éliminer les éventuelles légendes, notes de bas de page, etc... qui n'ont pas d'utilité.

En cas d'erreur, un bouton « RAZ » permet de restaurer la version pré-rognage.

Analyse des blocs sur plusieurs pages :

Il arrive qu'un bloc se prolonge sur plusieurs pages. Ceci est détecté en comparant les polices des derniers et premiers blocs des pages en question, leurs tailles de polices, ainsi que leurs coordonnées de début/fin. Les blocs ainsi repérés sont marqués.



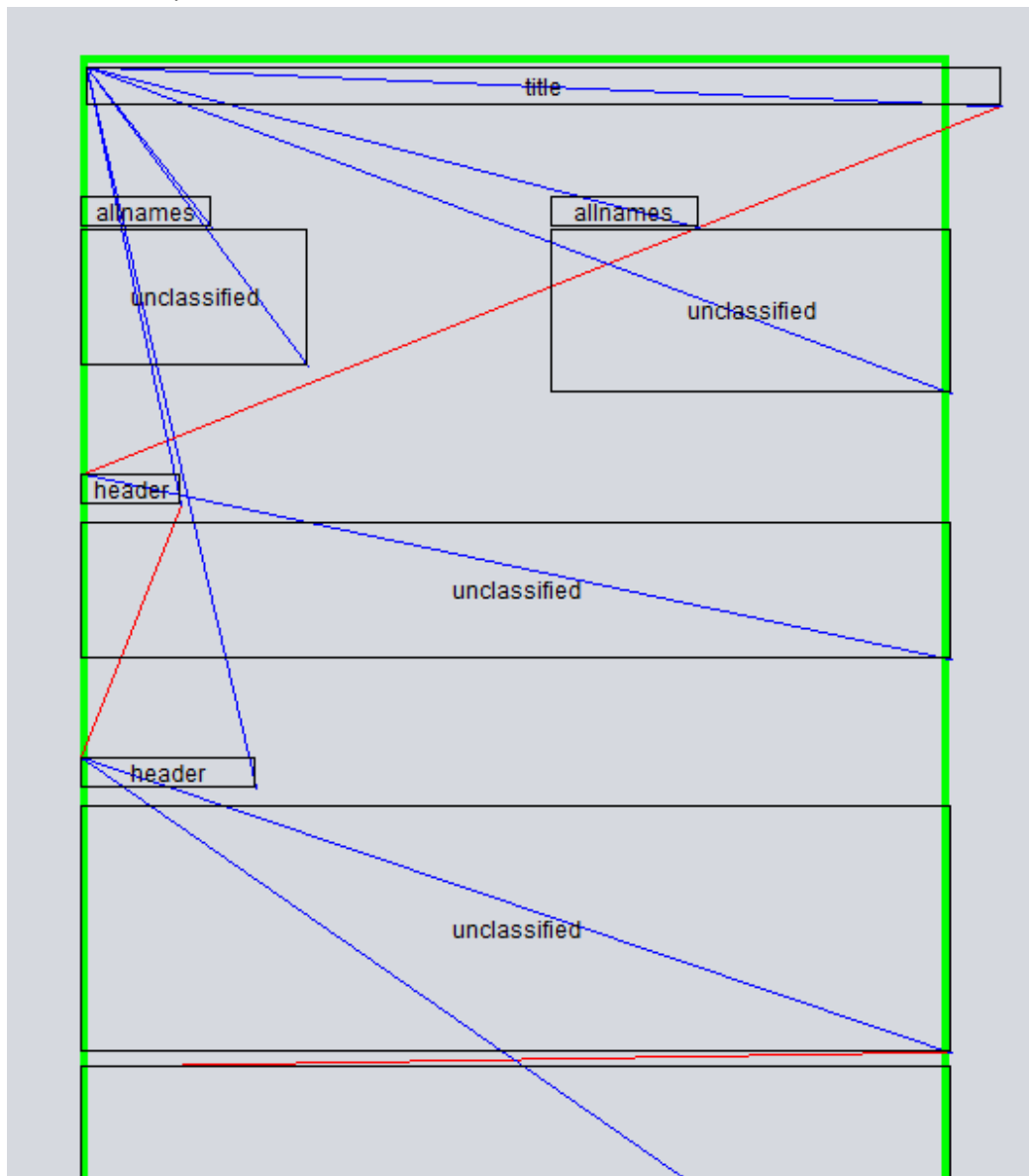
Note : Cette option est très inefficace si lancée avant d'avoir rogné la page. En effet, les premiers blocs des pages sont souvent les noms d'auteurs, des titres, des notes de bas de page...

Analyse hiérarchique :

Le but ici est de définir des relations père/fils, et frère/frère entre les blocs, afin d'extraire un arbre hiérarchique. Cette hiérarchisation est soumise à plusieurs règles, soumise à modification:

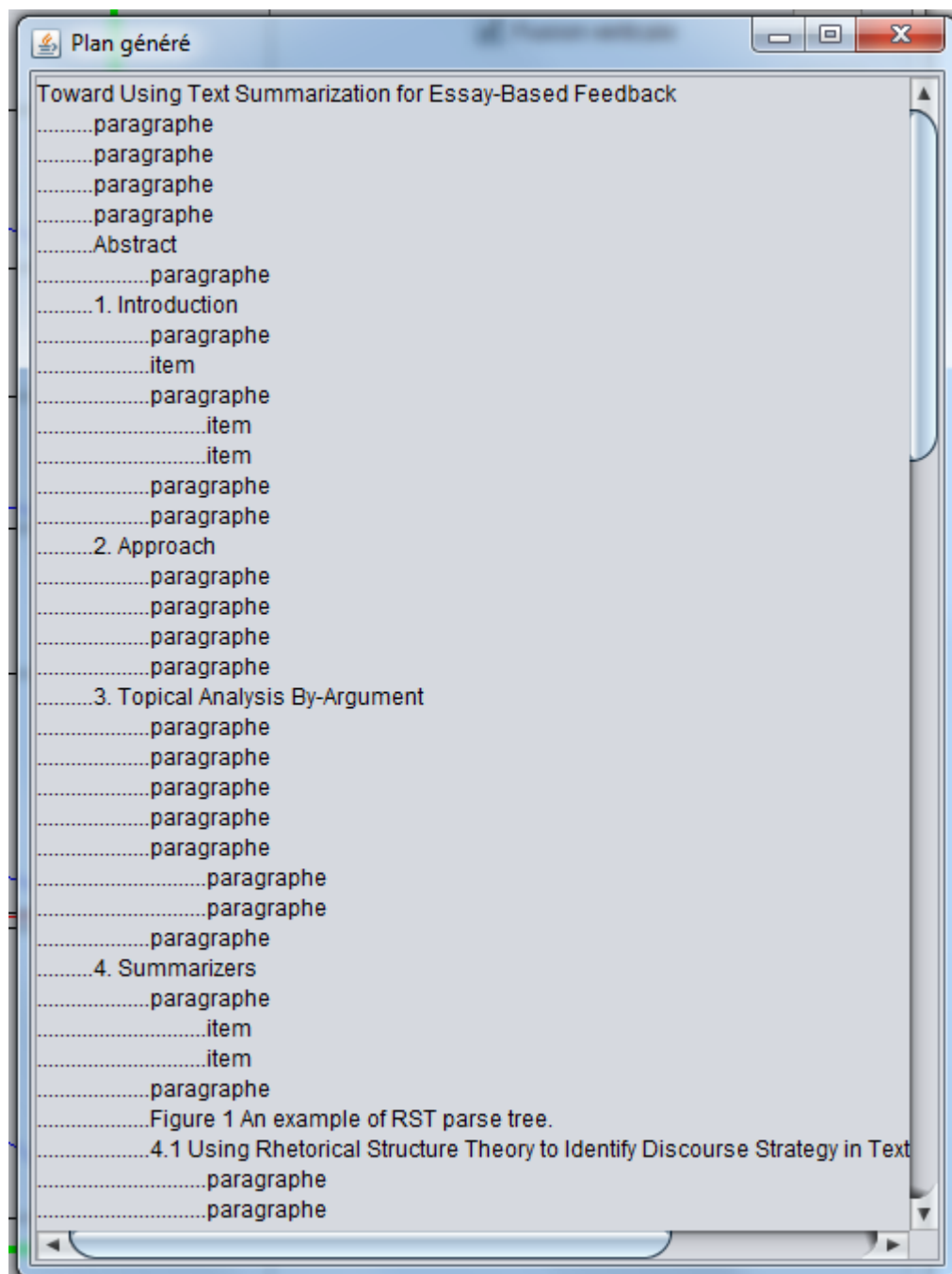
- Si un titre est au-dessus d'un non-titre, le titre est le père du non-titre.
- Deux titres ne sont frères que s'ils ont la même taille de police. Si ce n'est pas le cas, le titre le plus grand devient père de l'autre s'il est au dessus.
- Deux blocs non-titres ne sont frères que si leur coordonnée de départ est à peu près la même
- Si ce n'est pas le cas, et que le bloc du dessus commence avant celui du dessous, il devient le père de celui en dessous.
- Si ce n'est pas le cas, et que le bloc du dessus commence après celui du dessous, il devient le frère du père de celui au-dessus, si et seulement si celui du dessus avait déjà un père.
- Tout bloc qui n'a pas de père est automatiquement fils du premier bloc de la première page.
- L'analyse agit bloc par bloc, en comparant le prochain bloc non analysé au tout dernier bloc analysé. S'il n'est pas possible de se prononcer, l'analyse est relancée entre le même bloc non analysé et le père du tout dernier bloc analysé. La fonction remonte ainsi de père en père.

Voici un exemple :



L'affichage graphique est assez complexe à appréhender mais on s'y habitue vite. Les relations bleues représentent une relation père-fils, et les relations rouges une relation frère-frère. Ici, on voit que le premier bloc est père de tous les titres qui n'ont pas d'autres pères, comme le disent les règles. Les titres sont tous frères entre eux. Les deux blocs à la fin sont frères, car ils sont semblables.

Une fois les relations de hiérarchie établies, la personne peut choisir de générer un plan automatique du document.



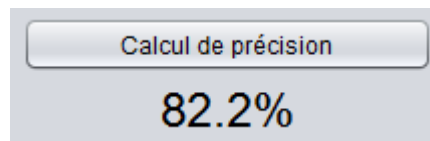
Annotation manuelle :

Si l'utilisateur le souhaite, il peut annoter manuellement son PDF afin d'évaluer la précision du programme. Pour ce faire, une interface spéciale est nécessaire : on y accède grâce au bouton « Annoter manuellement ».



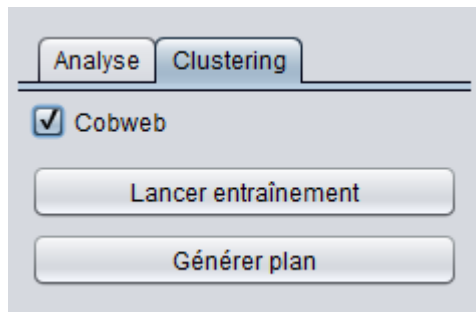
Sur cette interface, on distingue le PDF sous format d'image, deux boutons pour changer de page, une zone de texte pour assigner un type et un cadre rouge qui indique le bloc dont il est question.

Une fois cette annotation effectuée, l'utilisateur peut calculer une précision qui lui indiquera le pourcentage de types automatiquement assignés qui sont corrects, en se basant sur les types manuellement assignés.



Clustering :

Un outil de clustering a été ajouté au programme. Celui-ci permet de classifier les blocs suivant une approche purement algorithmique et non suivant des règles pré-établies. On accède aux options en question dans l'onglet « Clustering » du menu.



Pour l'instant seul le clusterer CobWeb est disponible. On lance l'entraînement de celui-ci via le bouton dédié, et on peut générer un plan par la suite. Le plan est un graphe reprenant tous les clusters par ordre hiérarchique.