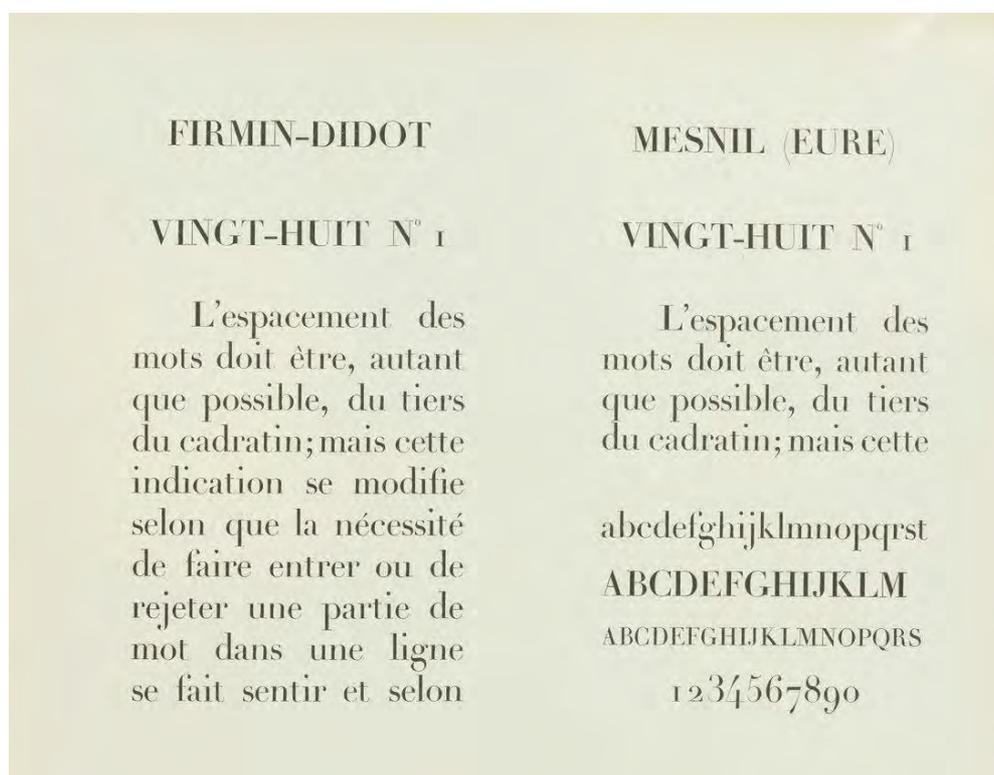


# Référentiel OCR



|                                       |                                       |
|---------------------------------------|---------------------------------------|
| Bibliothèque nationale de France      |                                       |
| direction des Services et des réseaux | Date :le mardi 28 avril 2015          |
| département de la Conservation        | Version :2                            |
| service Numérisation                  | Référence BnF :BnF-ADM-2014-062821-02 |

## TABLE DES MATIERES

|           |  |           |
|-----------|--|-----------|
| <b>1.</b> | <b>INTRODUCTION</b>  | <b>5</b>  |
| 1.1       | CONTEXTE   | 5         |
| 1.2       | OBJET  | 5         |
| 1.3       | DOMAINE D'APPLICATION  | 5         |
| <b>2.</b> | <b>DOCUMENTS APPLICABLES ET DE REFERENCE</b>                 | <b>6</b>  |
| <b>3.</b> | <b>GENERALITES</b>   | <b>7</b>  |
| 3.1       | OCR  | 7         |
| 3.1.1     | Documents éligibles à la conversion OCR                      | 8         |
| 3.2       | FORMAT ALTO  | 8         |
| 3.2.1     | Introduction   | 8         |
| 3.2.2     | Présentation du format ALTO BnF                              | 8         |
| 3.3       | NIVEAUX DE QUALITE   | 11        |
| 3.3.1     | Segmentation et structuration                                | 11        |
| 3.3.2     | Reconnaissance OCR   | 11        |
| <b>4.</b> | <b>RECONNAISSANCE DU TEXTE</b>                               | <b>13</b> |
| 4.1       | LANGUES ET ENCODAGE DE CARACTERES                            | 13        |
| 4.1.1     | Détection de la langue                                       | 13        |
| 4.1.2     | Traitement   | 14        |
| 4.1.3     | Taux qualité OCR   | 16        |
| 4.1.4     | Encodage   | 19        |
| 4.1.5     | Signes typographiques, caractères et symboles spéciaux, etc. | 19        |
| 4.2       | STYLES TYPOGRAPHIQUES  | 20        |
| 4.2.1     | Niveau de qualité  | 20        |
| 4.2.2     | Exposant, indice   | 20        |
| 4.2.3     | Titres   | 20        |
| 4.2.4     | Typographies mal reconnues par l'OCR                         | 20        |
| 4.3       | CESURES  | 22        |
| 4.3.1     | Répartition entre TextBlock                                  | 23        |
| <b>5.</b> | <b>SEGMENTATION ET STRUCTURATION</b>                         | <b>24</b> |
| 5.1       | DESCRIPTION DES PAGES  | 24        |
| 5.1.1     | Numérotation des pages                                       | 24        |

|        |  |           |
|--------|--|-----------|
| 5.1.2  | Qualité de numérisation des pages  | 25        |
| 5.1.3  | Qualité de l'océrisation des pages   | 25        |
| 5.1.4  | Pages vides  | 25        |
| 5.1.5  | Pages avec contenu en marges mais sans contenu principal                     | 26        |
| 5.1.6  | Pages de logo  | 26        |
| 5.2    | ORIENTATION DE LA PAGE   | <b>27</b> |
| 5.3    | STRUCTURATION DE LA PAGE   | <b>30</b> |
| 5.3.1  | PrintSpace   | 30        |
| 5.3.2  | XxxMargin  | 30        |
| 5.4    | BLOCS MANQUES  | <b>31</b> |
| 5.5    | ORDRE DE LECTURE ET ORDRE DES SEGMENTS                                       | <b>32</b> |
| 5.5.1  | Mise en page en colonnes   | 32        |
| 5.5.2  | Mise en page en colonnes avec des éléments centraux non textuels             | 34        |
| 5.5.3  | Mise en page en colonnes avec des éléments centraux textuels                 | 35        |
| 5.5.4  | Corps du texte et notes séparés par un trait                                 | 38        |
| 5.6    | TEXTE  | <b>38</b> |
| 5.6.1  | Paragraphes  | 38        |
| 5.6.2  | Titres   | 39        |
| 5.6.3  | Tableaux   | 39        |
| 5.6.4  | Encadrés   | 43        |
| 5.6.5  | Notes de bas de page   | 44        |
| 5.6.6  | Illustrations avec habillage de texte traversant                             | 45        |
| 5.6.7  | Publicités et catalogues d'éditeur   | 45        |
| 5.6.8  | Texte sous tampon  | 47        |
| 5.6.9  | Texte illisible  | 48        |
| 5.6.10 | Texte manqué   | 49        |
| 5.6.11 | Texte en OCR brut  | 49        |
| 5.7    | ILLUSTRATIONS  | <b>49</b> |
| 5.7.1  | Typage « illustration »  | 51        |
| 5.7.2  | Imbrication de blocs Illustration et d'autres blocs, notamment des TextBlock | 52        |
| 5.7.3  | Formules chimiques, mathématiques  | 54        |
| 5.7.4  | Partitions   | 54        |
| 5.7.5  | Cartes   | 55        |
| 5.7.6  | Alphabets non latins   | 55        |
| 5.7.7  | Illustration en écriture manuscrite  | 55        |
| 5.8    | ELEMENTS GRAPHIQUES  | <b>57</b> |
| 5.8.1  | Décorations et ornements   | 57        |
| 5.8.2  | Tampons  | 57        |

|           |  |           |
|-----------|--|-----------|
| 5.8.3     | Lettrines (lettres ornées)   | 58        |
| 5.8.4     | Traits de séparation   | 58        |
| 5.8.5     | Ecriture manuscrite  | 63        |
| 5.8.6     | Imbrication de blocs GraphicalElement et d'autres blocs, notamment des TextBlock | 63        |
| 5.9       | BLOCS COMPOSES   | <b>63</b> |
| 5.9.1     | Texte au sein des illustrations ou des éléments graphiques                       | 63        |
| 5.9.2     | Imbrication d'illustrations ou d'éléments graphiques et de texte                 | 64        |
| 5.9.3     | Ordre de lecture entre texte et illustrations ou éléments graphiques             | 65        |
| 5.10      | TABLEAU RECAPITULATIF DE LA STRUCTURATION ALTO                                   | <b>65</b> |
| 5.10.1    | Etiquetage des éléments  | 66        |
| <b>6.</b> | <b>QUALITE DE LA RECONNAISSANCE OCR</b>  | <b>71</b> |
| 6.1       | QUALITE DE LA TRANSCRIPTION DU TEXTE   | <b>71</b> |
| 6.1.1     | Correction ciblée  | 71        |
| 6.2       | QUALITE DE LA SEGMENTATION   | <b>72</b> |
| 6.3       | DEQUALIFICATION DE CONTENUS DANS UN DOCUMENT                                     | <b>72</b> |
| 6.3.1     | Déqualification par types de contenu   | 72        |
| 6.3.2     | Déqualification des mots ou blocs illisibles                                     | 73        |
| 6.3.3     | Limites au principe de déqualification   | 73        |
| 6.4       | DEQUALIFICATION DU TAUX QUALITE SUR UN DOCUMENT                                  | <b>73</b> |
| 6.5       | DEQUALIFICATION OU REFUS DE DOCUMENTS  | <b>74</b> |
| <b>7.</b> | <b>CONTROLE DE LA QUALITE</b>  | <b>75</b> |
| 7.1       | CONTROLE AUTOMATIQUE ALTO  | <b>75</b> |
| 7.2       | CONTROLE PAR ECHANTILLONNAGE VISUEL  | <b>76</b> |
| 7.2.1     | Qualité de la reconnaissance du texte  | 76        |
| 7.2.2     | Qualité de la segmentation/structuration   | 77        |
| 7.2.3     | Détail des métriques qualité   | 79        |

# 1. INTRODUCTION

---

## 1.1 Contexte

La Bibliothèque nationale de France a lancé divers programmes concourant à la constitution d'une bibliothèque numérique. Ces programmes s'appuient notamment sur des marchés de dématérialisation des collections de la BnF et de bibliothèques françaises partenaires.

En sus de la numérisation proprement dite des documents, un certain nombre d'autres prestations de dématérialisation sont demandées, dont la reconnaissance et la conversion des contenus textuels des documents numérisés.

En effet, la BnF désire également donner accès au contenu des documents grâce à la recherche à partir du texte des pages. La mise en place de ce procédé implique donc la conversion en mode texte de l'intégralité du contenu des pages afin de permettre la recherche plein texte, quelle que soit la partie consultée, puis l'affichage des images du document correspondant, avec possibilité d'accéder aux données en mode texte pour faire des sélections, des copies, des impressions.

Cette conversion en mode texte s'appuie principalement sur des techniques de reconnaissance optique de caractères (OCR, *optical character recognition*).

Ce document est organisé en plusieurs parties :

- Un chapitre « Généralités » qui expose les principes généraux de la reconnaissance optique de caractères.
- Un chapitre « OCR » qui traite de l'extraction des contenus.
- Un chapitre « Segmentation et structuration » qui traite de l'extraction de la structure du document.
- Un chapitre « Qualité » qui décrit les principes et méthodes de mesure de la qualité.

## 1.2 Objet

Le présent référentiel définit les caractéristiques attendues pour le traitement de reconnaissance optique de caractères appliqué aux documents des départements de Bibliothèque nationale de France et des bibliothèques partenaires. Il détaille les caractéristiques techniques des fichiers, les modalités de contrôle, etc.

## 1.3 Domaine d'application

Le présent référentiel s'applique aux prestations de numérisation d'ouvrages commandées par la Bibliothèque nationale de France, sur des marchés de numérisation ou de réfection, ainsi qu'aux éventuelles productions internes.

## 2. DOCUMENTS APPLICABLES ET DE REFERENCE

---

|   |   |
|---|---|
| Schéma ALTO LoC                                   | <a href="http://www.loc.gov/standards/alto (version 3.0)">http://www.loc.gov/standards/alto (version 3.0)</a>   |
| Schéma ALTO BnF                                   | <a href="http://bibnum.bnf.fr/schema/alto_bnf-v2_0.xsd">http://bibnum.bnf.fr/schema/alto_bnf-v2_0.xsd</a>   |
| BnF : Conversion en mode texte                    | <a href="http://www.bnf.fr/fr/professionnels/num_conversion_texte/s.num_conversion_texte_ocr.html">http://www.bnf.fr/fr/professionnels/num_conversion_texte/s.num_conversion_texte_ocr.html</a>                     |
| Référentiels BnF                                  | <a href="http://www.bnf.fr/fr/professionnels/numerisation_boite_outils/a.numerisation_referentiels_bnf.html">http://www.bnf.fr/fr/professionnels/numerisation_boite_outils/a.numerisation_referentiels_bnf.html</a> |
| Référentiel d'enrichissement des métadonnées      | version 3 / BnF-ADM-2013-079381-04  |
| Référentiel d'enrichissement des métadonnées METS | version 1 / BnF-ADM-2013-117422-01  |
| Référentiel de livraison de document numérique    | version 3 / BnF-ADM-2013-077351-03  |
| Référentiel Tables                                | version 1 / BnF-ADM-2014-062888-01  |

### 3. GENERALITES

---

#### 3.1 OCR

La reconnaissance optique de caractères désigne les procédés informatiques visant à extraire le texte présent dans l'image d'un texte imprimé. Un système OCR part donc de l'image numérique réalisée par un scanner optique ou une caméra numérique d'une page (document imprimé, feuillet dactylographié, documents transparents, etc.), et produit en sortie un fichier texte en divers formats.

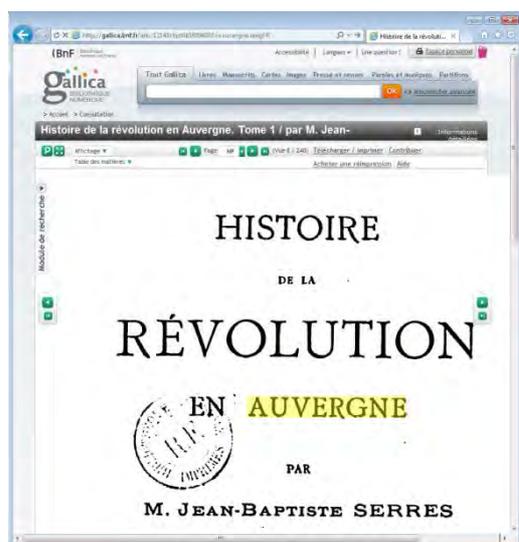
Ces systèmes OCR opèrent en plusieurs étapes :

1. *Pré-analyse de l'image*, visant à améliorer la qualité de l'image en vue de faciliter la reconnaissance des caractères (redressement d'images inclinées, des corrections de contraste, binarisation de l'image).
2. *Segmentation des contenus*, permettant d'isoler dans l'image les différentes composantes (illustrations, blocs de texte, marges, etc.).
3. *Reconnaissance des caractères* : le texte contenu dans les blocs isolés à l'étape précédente est analysé et reconnu.



La segmentation physique de la page (étape 2) permet de mettre en correspondance le texte issu de l'OCR avec son emplacement dans l'image de la page. Cette correspondance est utilisée lors de la consultation des documents de la bibliothèque numérique Gallica, en permettant la mise en valeur, sur l'image du document, du ou des mots-clés recherchés.

Elle est également utilisée lors de la génération de documents PDF restituant le document numérisé (en mode image) mais aussi son contenu textuel (afin de rendre possible la recherche en texte intégral).



*Recherche en texte intégral et mise en surbrillance du critère de recherche dans Gallica*

---

### 3.1.1 Documents éligibles à la conversion OCR

Les collections de la BnF sont constituées d'une grande variété de typologies d'ouvrages. Du fait des limitations intrinsèques des procédés informatiques utilisés pour la conversion OCR, seuls un sous-ensemble des documents patrimoniaux sont éligibles à ce traitement (voir section 3.3).

La commande d'une conversion en mode texte pour un document donné a généralement lieu lors de l'état conjoint, pendant lequel les représentants de la BnF et le prestataire précisent la commande en fonction de la nature du document.

## 3.2 Format ALTO

### 3.2.1 Introduction

ALTO (*Analysed Layout and Text Object*) est un standard XML permettant de rendre compte de la mise en page physique et de la structure logique d'un texte reconnu par un système OCR. Ce format est issu du projet européen METAe1 et il est actuellement maintenu par un comité éditorial hébergé par la Bibliothèque du Congrès (<http://www.loc.gov/standards/alto/>).

ALTO est très utilisé pour la conversion en mode texte de documents patrimoniaux, en France et à l'étranger. Il est bien adapté à la conservation à long terme des données issues de la conversion et il permet une réutilisation ultérieure du mode texte, dans la mesure où il contient pour chaque mot et bloc de texte ses coordonnées dans la page, le taux de confiance de reconnaissance, éventuellement des éléments de forme (styles de caractère, polices).

La Bibliothèque nationale de France a introduit des restrictions dans le format ALTO originel et utilise donc une variante pour sa production. La nature de ces changements est documentée dans le schéma lui-même.



#### Pourquoi le format ALTO et non le PDF ?

La BnF a préféré le format ALTO au format PDF comme support de l'action de dématérialisation en mode texte des documents patrimoniaux car il s'agit :

- d'un format XML, avec les qualités intrinsèques du monde XML : universalité, facilité de création, d'édition et d'archivage, compacité, etc.
- d'un format conçu expressément pour l'usage attendu (numérisation patrimoniale en modes image et texte) alors que le PDF répond à l'origine aux besoins liés à l'impression de documents numériques. Il est possible de générer le PDF à partir des images et des fichiers ALTO alors que l'inverse n'est pas possible.

### 3.2.2 Présentation du format ALTO BnF

Le format est organisé en cinq grandes composantes :

- En-tête XML
- Description de la page et de son traitement
- Styles détectés dans la page
- Etiquettes présentes dans la page

- Segmentation de la page

### *Entête XML*

L'élément <alto> donne les informations suivantes :

- Espaces de noms utilisés au sein du schéma ALTO BnF : xlink, ALTO LoC
- Version du schéma ALTO BnF : l'attribut SCHEMAVERSION donne la version du schéma utilisé, par exemple **SCHEMAVERSION="alto\_bnf-v2\_0"**

### *Description de la page et de son traitement*

L'élément <Description> permet de fournir les informations suivantes :

- unité de mesure utilisée (<MeasurementUnit>) : pixel
- information sur l'image source océrisée :
  - nom du fichier image (<fileName>) : le motif imposé est `\d{8}.(TIF|tif|JPG|jpg|jp2|JP2)`
  - identifiant du fichier image dans la bibliothèque numérique, si cet identifiant est connu (dans le cas d'un traitement complémentaire) : `<fileIdentifier>ark:/12148/bpt6k75043976/f1</fileIdentifier>`
  - identifiant de production du document numérique : les deux catégories d'identifiants producteurs à utiliser sont « NUM » et « IFN » : `<documentIdentifier>documentIdentifierLocation="NUM">7504397</documentIdentifier>`
- procédé de traitement OCR (<OCRProcessing>, <ocrProcessingStep>) :
  - date du traitement (<processingDateTime>) : format ISO AAAA-MM-JJ
  - information de production (<processingStepDescription>) : nombre de caractères, de chaînes (cf. section 7.2.3), types de traitement appliqués (par exemple correction des courbures de ligne), etc.
  - paramètres du traitement (<processingStepSettings>) : moteur OCR utilisé, version, réglages spécifiques, etc.

### *Styles détectés dans la page*

L'élément <Styles> permet de décrire les styles détectés dans la page (cf. section 4.2).

### *Étiquettes présentes dans la page*

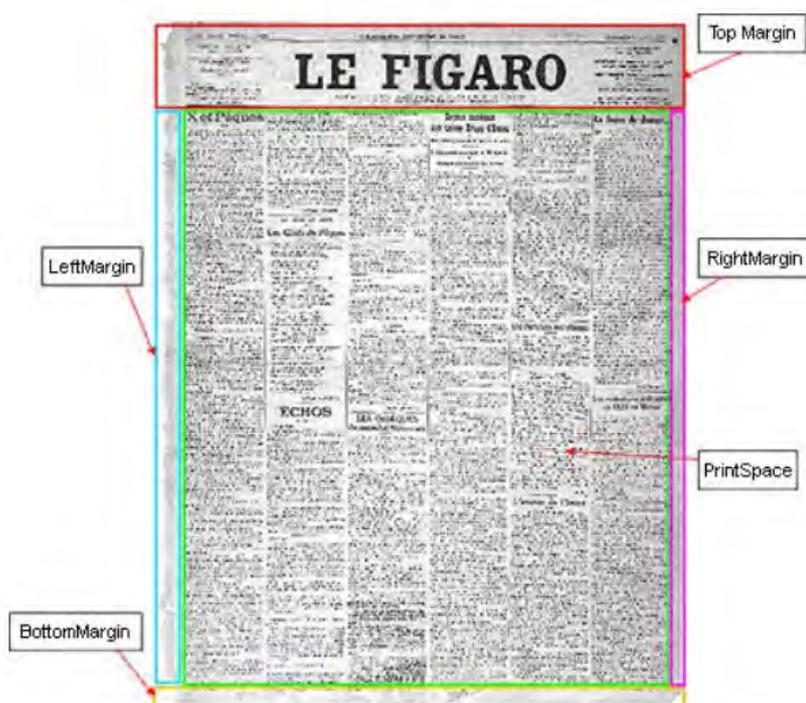
L'élément <Tags> permet de décrire les étiquettes présentes dans la page (cf. section 5.10.1). Ces étiquettes caractérisent divers types de contenus : entités nommées, titres, etc.

### *Segmentation de la page*

Le format ALTO permet de décrire la segmentation d'une page (<Layout>) en ses différentes composantes. L'élément <page> peut ainsi contenir cinq composantes :

- TopMargin : désigne la zone supérieure de la page du bord gauche au bord droit hors zone de texte. Quand c'est possible, il s'agit de la zone contenant le titre, l'ours, etc.

- BottomMargin : désigne la zone inférieure de la page du bord gauche au bord droit hors zone de texte.
- LeftMargin : désigne la zone gauche de la page hors zone supérieure, zone inférieure et zone de texte.
- RightMargin : désigne la zone droite de la page hors zone supérieure, zone inférieure et zone de texte.
- PrintSpace : désigne la zone de texte. Cet élément est obligatoire. Il contient au moins un élément de type bloc.



*Exemple de découpage d'une page de presse*

Dès que l'un de ces éléments contient une information (texte, illustration etc.), cette information est décrite dans un ou plusieurs éléments de type bloc.

Les blocs peuvent être de quatre types différents :

- TextBlock : désigne le bloc de texte. Cet élément est utilisé pour regrouper les lignes de textes en un ensemble cohérent.
- Illustration : désigne une image ou une figure.
- GraphicalElement : désigne un élément graphique autre qu'une image ou une figure. Il peut être utilisé pour décrire un élément de séparation intertextuel ou un élément textuel non reconnu en tant que tel par l'OCR.
- ComposedBlock : est utilisé pour permettre l'imbrication d'éléments bloc.

A l'intérieur d'un TextBlock, l'élément Line décrit les lignes de texte et l'élément String rassemble les caractères en mots.

Les coordonnées des éléments sont définies à partir du point de repère le plus en haut à gauche de la page. Ainsi, chaque bloc, ligne ou chaîne de caractères reconnus est identifié dans l'ordre de présentation de l'original.

ALTO permet également de décrire des formes géométriques (cercle, polygone, ellipse), de gérer les césures... Les objets non textuels ont également leurs propres découpage et coordonnées.

#### *Identification des mots et chaînes de caractères*

Chaque chaîne de caractères composant un mot ou une partie de mot césuré (String) est identifiée avec les informations suivantes :

- contenu : le mot reconnu par le système OCR et/ou corrigé manuellement selon le niveau de qualité demandé ;
- *wc (word confidence)* : note de confiance de la reconnaissance de chaque mot, notée de 0 à 1 ;
- *cc (character confidence)* : note de confiance de la reconnaissance de chaque caractère du mot. Cette note est composée d'une liste de notes de 0 à 9, une note pour chaque caractère ;
- *wd (word dictionary, optionnel)* : appartenance ou non du mot à un dictionnaire.

### **3.3 Niveaux de qualité**

Les niveaux de qualité attendus concernent tant la qualité de la segmentation que la qualité de la reconnaissance OCR.

#### **3.3.1 Segmentation et structuration**

La qualité de la segmentation et de la structuration des contenus du document d'origine lors de leur transcription au format ALTO concerne en particulier :

- l'ordre de lecture,
- le typage des blocs,
- le chevauchement des blocs.

La qualité de la segmentation et de la structuration est détaillée aux sections 7.1 et 7.2.1.

#### **3.3.2 Reconnaissance OCR**

Le résultat de la reconnaissance OCR est très variable, en fonction de la nature bibliographique et physique des documents océrisés :

- lisibilité du document : les défauts d'impression, le vieillissement du papier, les problèmes de migration d'encre ou de courbure de page, etc., influent sur la qualité de l'image numérisée et donc sur l'aptitude des systèmes OCR à en extraire le texte ;
- critères bibliographiques : la qualité de reconnaissance des systèmes OCR est particulièrement sensible à la nature des contenus, en termes de langue et d'alphabet notamment ;
- genres documentaires : les systèmes OCR s'appuient sur des dictionnaires pour affiner leur reconnaissance des mots. Idéalement, il faudrait donc disposer de dictionnaires contemporains de la date d'édition et adaptés au contenu de chaque ouvrage (littérature, sciences, philosophie, etc.).



Différents taux sont attendus selon les marchés (voir section 4.1.3) :

- OCR haute qualité : par exemple 99,9 %
- OCR taux qualité garantie : par exemple 98,5 %
- OCR brut : pour les ouvrages datant de 1750 ou antérieur, ou pour les ouvrages sans date.

Ces taux sont précisés dans le contexte de chaque marché de numérisation.

Cette qualité de reconnaissance OCR est évaluée à l'aide d'un taux de confiance (wc, voir section 3.2.2) et non d'un taux effectif. En effet, pour connaître le taux de qualité effectif pour un document océrisé, il faudrait avoir connaissance de sa vérité terrain (le texte exact du document) afin de la comparer avec le texte produit par le logiciel OCR, ce qui est bien sûr impossible dans le cadre d'une numérisation de masse. Ce taux de confiance est fourni par les moteurs OCR, pour chaque mot traité par le logiciel.

Les taux de qualité OCR BnF sont calculés sur la base du mot, en moyennant les taux de confiance des mots présents dans un document.



Le taux de qualité d'un document n'est pas obtenu en moyennant les taux de qualité de chaque page, mais bien en moyennant les taux de confiance de tous les mots du document.

Un document peut avoir des parties dont le taux de reconnaissance est supérieur ou inférieur au taux qualité admissible, mais la moyenne doit correspondre à la qualité exigée.

La qualité de la reconnaissance OCR est détaillée à la section 6.1.

Le contrôle de la qualité de la reconnaissance OCR est détaillé à la section 7.2.2.

## 4. RECONNAISSANCE DU TEXTE

---

Ce chapitre présente des informations sur l'encodage du texte et le traitement des différentes polices et langues ainsi que sur la gestion des césures.

### 4.1 Langues et encodage de caractères

#### 4.1.1 Détection de la langue

Le moteur OCR permet de définir la langue et les caractères pouvant faire partie du contenu (soit à l'ouvrage, soit à la page, soit au bloc texte). Ceci permet d'associer le dictionnaire de la langue du texte et de déterminer la fiabilité de la reconnaissance.

Tous les ouvrages sont traités par défaut comme étant en langue française, sauf détection par le moteur OCR d'autres langues.

La langue de chaque bloc de texte sera consignée dans l'attribut LANG du schéma ALTO, selon la norme ISO 639-2 :

LANG="fr"



Si deux blocs de langues différentes se suivent, leur langue sera détectée et consignée.

A MARIE.

Montagne aux doux parfums, échelle radieuse,  
AMOUR, LOUANGE ET GLOIRE.

Maria mons in quo beneplacitum est Deo habitare.

Marie est cette montagne où Dieu prend plaisir d'habiter.      ST. ATHANASE.

Maria, scala quant Jacob vidit ad cælos pertingentem.

Marie est cette échelle que vit Jacob et qui allait jusqu'au ciel.      Idem.

❦

Au ciel, si je veux m'élancer,  
Je sens mon impuissance ;  
Mon cœur veut en vain s'efforcer,  
Je perds toute espérance :  
Venez, Marie, à mon secours,  
Car c'est à vous que j'ai recours.

*Texte en latin dans une page majoritairement en français*

---

Si deux langues sont présentes dans le même bloc, la langue majoritaire sera consignée au niveau du bloc et la langue minoritaire au niveau des mots.

|   |   |
|---|---|
| Ἐγκώμιον (le genre), 232-240.                 | Ἐπιπέστωσις, 26.                            |
| Ἐκλογή, qualité du style, 422.                | Ἐπιτροπή, 507, n. 2.                        |
| Ἐιδωλοποιία, 503, n. 2.                       | Épisodes dans l'éloge, 308.                 |
| Ἐἶδη (les), 123.                              | Ἐπιστολικόν (le genre), 236-244.            |
| Ἐἶκος, p. 115.                                | Ἐπιτάφιος λόγος, 241.                       |
| Ἐκφώνησις, 507, n. 3.                         | Épithètes (les), 476.                       |
| Ἐλεγγος, 41-151.                              | Ἐπίτροχον ou ἐπιδρομή, 464.                 |
| Élégance (l'), 487.                           | Ἐπίθετα (les mots), 21.                     |
| Ἐλεσοί (les), 23.                             | Ἐπίθετοι (les qualités), 430.               |
| Ἐλληνικοὶ λόγοι, 33.                          | Ἐπόρωσις, 21.                               |
| Ἐλληγίζειν (τό), 422.                         | Équité (de l'), 342.                        |
| Élocution (l'), 96-413.                       | Ἐκκληξίς, 466.                              |
| Éloges (les), 241.                            | Ἐρμήνευται, 413.                            |
| Éloquence appelée philosophie, 32.            | Esprit (l'), 487.                           |
| Éloquence (définition de l'), 75-88           | Ἐσχηματισμένοι λόγοι, 255.                  |
| — (Rapports de l') et de la dialectique), 83. | Esthétique (l'élément) dans l'éloque), 352. |

*Blocs en français, en grec et mixtes*

#### 4.1.2 Traitement

Les langues romanes et les langues écrites avec l'alphabet latin (y compris le latin) seront traitées en OCR, ainsi que le grec.

Pour le français, le prestataire décrira son aptitude à mettre en œuvre des dictionnaires d'orthographe adapté à la date de publication des documents : ancien français, moyen français, français classique, français moderne. Il exposera en outre sa capacité à traiter les formes anciennes de la lettre *s* minuscule : *s* long, (l), eszett (ß).

Pour le traitement des langues suivantes, le prestataire décrira son aptitude à mettre en œuvre des dictionnaires d'orthographe : allemand, anglais, italien, espagnol, grec, grec ancien, latin, portugais. Le taux qualité attendu pour ces langues sera l'OCR brut.

Les langues écrites avec l'alphabet grec (autre que le grec) et les langues non romanes écrites avec l'alphabet latin (albanais, croate, estonien, finnois, hongrois, lithuanien, roumain, slovaque, slovène, norvégien, polonais, tchèque, turc moderne, vietnamien en écriture latinisé, etc.), seront également traités en OCR, mais le taux qualité attendu sera l'OCR brut, sauf demande particulière de la BnF et selon les propositions techniques du prestataire.

Les autres alphabets et types d'écriture (cyrillique, langues asiatiques, arabe, hébreu, etc.) ne seront pas traités en OCR. Les blocs de texte où ces alphabets ou systèmes d'idéogrammes sont présents seront décrits sous la forme de blocs image. Ces blocs seront des éléments Illustration doté d'une étiquette LayoutTag LABEL="nonLatinScript" (cf. section 5.10).



### EXEMPLE

|                           |                           |
|---------------------------|---------------------------|
| أفريقي ١١٢, 17            | بركس بنارس ١٩٠, 5         |
| أفغور شاه ١١٣, 2 — ١١٩, 8 | برخوشيا v. برخوشيا ١٣١, 5 |
| أكسيرخس ١٩٠, 7            | بركومنس ١٩٠, 1            |
| أكسيوتس ١٩٠, 1            | بلاسوس ١٩٤, 9             |
| أنتي ثودي ١١٩, 11         | بلدة الثعلب ١٣١, 17       |
| أنهاء التجارة ١٣٨, 3. 8   | بليلج ١٣٠, 4              |
| امتلاء ١٧١, 9 — ١٧٣ — ١٧٥ | بليناس ١٨٤, 18            |
| املج ١٨٣, 4               | بهارات ١٨٩, 18            |
| أحمركانيك ١٣٧, 22         | بورنشيا ١٦٩, 5            |
| أنوشيروان ١٣٩, 11         | بيت ١٣٨, 1 ff.            |
| الانيسلان ١٣٥, 18         | تابع النجم ١٣٤٢, 15       |
| أخليلج ١٨٣, 4             | تأسيس ١٣٤٠, 22            |
| أونرساوس (?) ١٣٨, 2       | الثغثي ١٣٦٢, 18 — ١٣٥١, 8 |

Bloc de texte à traiter en mode image

Les blocs de texte où sont présents des langues en alphabet latin et non latin (cas des dictionnaires de langues par exemple) seront traités en OCR afin de transcrire les mots en alphabet latin ; la langue des mots en alphabet non latin sera consignée au niveau des mots.



### EXEMPLE

| INDEX                                  |                | Pages  |
|--|----------------|--|
| Ái 脉.....                              | 29             | Chiêu 昭 (fils de Gia-Long) 43                                |
| An 安 (fils de Nguyễn Phúc-Nguyễn)..... | 19             | Chương 曄 appelé aussi Trà 茶..... 32                          |
| An 安 (fils de Minh-Mạng).....          | 44             | Chường 種..... 10   |
| An-làng 安陵.....                        | 13, 14, 67, 68 | Công-Thượng-Vương... 3 ff.                                   |
| Anh 洪.....                             | 19             | Cơ-thánh 基聖..... 9, 10                                       |
| Bãng 版.....                            | 34             | Cự 矩..... 42   |
| Biện 鼻.....                            | 15, 65         | Diễn 演 (fils de Nguyễn Hoàng)..... 17                        |
| Bình 柄.....                            | 28             | Diễn 演 appelé aussi Hán 漢 (fils de Nguyễn Phúc-Tân)..... 21  |
| Bình 柄 appelé aussi Úc 旭               | 41             | Diệu 逸..... 22   |
| Bình 平.....                            | 37             | Diệu 曜..... 35   |
| Bồi-làng 倍陵.....                       | 15             | Du 澍 appelé aussi Nghiêm 驥 (fils de Nguyễn Phúc-Chú)..... 31 |
| Bữu 寶.....                             | 33             | Duán 駒..... 42   |
| Bữu-Côn 寶峴.....                        | 66             | Dực 昱 appelé aussi Bữu 寶..... 33                             |
| Bữu-Cương 寶岡.....                      | 65             |  |
| Bữu-Hào 寶濤.....                        | 67             |  |
| Bữu-Lân 寶麟.....                        | 39, 66         |  |
| Bữu-Liêm 寶廉.....                       | 66             |  |
| Bữu-Lợi 寶麟.....                        | 67             |  |

Blocs de texte mixte (français, vietnamien latinisé et écriture en sinogrammes)

### 4.1.3 Taux qualité OCR

#### *Taux OCR brut*

Le taux OCR brut est celui obtenu en sortie du moteur OCR. Il est calculé par le moteur OCR.

Ce taux est inscrit dans le manifeste du document numérique et dans le fichier ALTO chacune de ses pages (cf. section 7.2.3).

Ce taux peut s'appliquer à tout le document ou à des portions de document (dans le cas de zones déqualifiées, cf. section 6.3).

#### *Taux OCR qualité garantie ou taux OCR haute qualité*

Ces taux sont ceux attendus pour une prestation OCR avec des exigences particulières portant sur la qualité de la transcription du texte. Ces prestations impliquent une montée en qualité manuelle du texte pour les documents dont l'OCR brut est inférieur au taux qualité ciblé.

#### Périmètre

Tous les blocs de texte sont comptabilisés dans l'évaluation du taux OCR qualité garantie du document, exceptés :

- les blocs en langue autre que le français,
- les blocs déqualifiés (voir section 6.3)

#### Jeu de caractères

Les chiffres ainsi que les caractères non alphanumériques mis en évidence dans le tableau suivant ne sont pas comptabilisés dans l'évaluation du taux OCR qualité garantie. Ils constituent le « jeu de caractères hors garantie ».

|           | x0 | x1 | x2 | x3 | x4 | x5 | x6 | x7 | x8 | x9 | <u>xA</u> | <u>xB</u> | <u>xC</u> | <u>xD</u> | <u>xE</u> | <u>xF</u> |
|-----------|----|----|----|----|----|----|----|----|----|----|-----------|-----------|-----------|-----------|-----------|-----------|
| 2x        |    | !  | "  | #  | \$ | %  | &  | '  | (  | )  | *         | +         | ,         | -         | .         | /         |
| 3x        |    |    |    |    |    |    |    |    |    |    | :         | ;         | <         | =         | >         | ?         |
| 4x        | @  |    |    |    |    |    |    |    |    |    |           |           |           |           |           |           |
| 5x        |    |    |    |    |    |    |    |    |    |    | [         | \         | ]         | ^         | _         |           |
| 7x        |    |    |    |    |    |    |    |    |    |    | {         |           | }         | ~         |           |           |
| <u>Ax</u> |    | ı  | ç  | £  | ¤  | ¥  | ı  | §  | "  | ©  |           | «         | ¬         | •         | -         |           |
| <u>Bx</u> | °  | ±  |    |    | ´  | µ  | ¶  | •  | ¸  |    | »         |           |           |           |           | ¿         |
| Fx        |    |    |    |    |    |    |    | ÷  |    |    |           |           |           |           |           |           |

De plus, ne seront pas comptabilisées comme erreur :

- les erreurs concernant la reconnaissance des caractères accentués : substitutions ou omissions de diacritiques (« e » pour « é » par exemple) ;

- les erreurs concernant la casse : substitutions ou omissions de casse : « A » pour « a »
- les erreurs de reconnaissance sur les chiffres arabes : « 10 » pour « 16 ».

Par contre, sont considérées comme des erreurs :

- un caractère alphabétique transcrit erronément en un caractère hors garantie : « avlon » pour « avion »
- un caractère hors garantie collé à un mot transcrit erronément en un caractère alphabétique : « avioni » pour « avion; »
- un caractère hors garantie intramot non transcrit : « nétant » pour « n'étant »
- un caractère hors garantie intramot transcrit en un caractère alphabétique : « niétant » pour « n'étant »
- l'insertion parasite d'un caractère hors garantie dans un mot : « situe'e » pour « située »

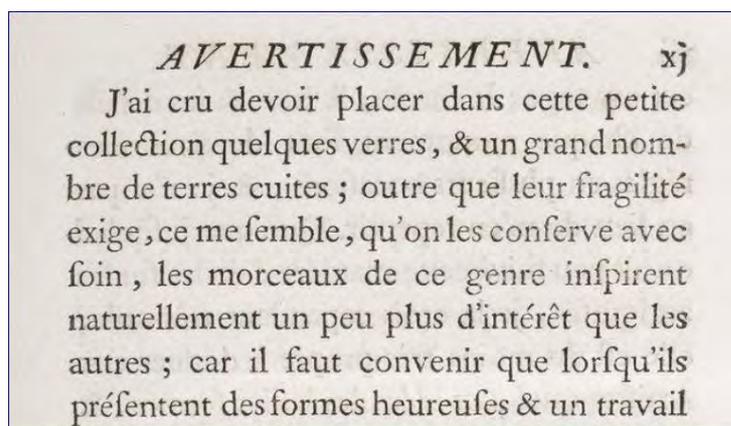


Une charte précisera ces règles en début de prestation OCR.

#### Typographie ancienne

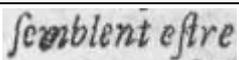
Les ouvrages composés avec des caractères d'imprimerie qui ne sont plus (ou peu) présents dans les imprimés contemporains relèvent de la catégorie « typographie ancienne » : s long (l), eszett (ß), ligatures, etc.

Des taux OCR spécifiques peuvent s'appliquer.



<http://gallica.bnf.fr/ark:/12148/btv1b8626613n>

Une translittération est exigée pour certains caractères, notamment les formes anciennes de la lettre s minuscule : s long, (l), eszett (ß). Cette translittération fait partie du périmètre du taux OCR qualité garantie.

| Glyphe | Exemple   | Translittération |
|--------|---|------------------|
| f/s    |  | semblent estre   |
| ff/ss  | <b>Palissades</b>   | Palissades       |



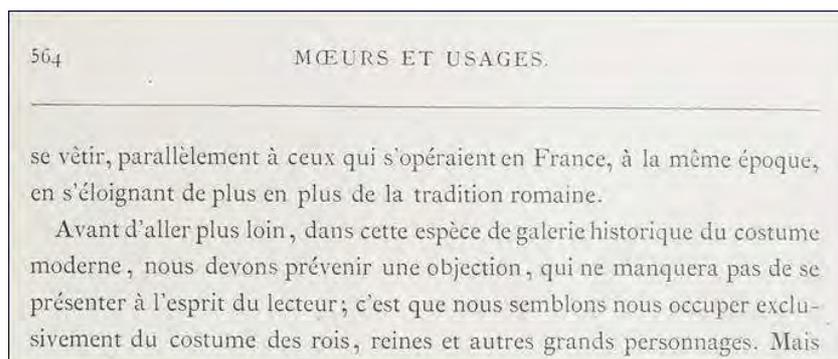
## NOTES

Il n'est pas demandé de dissimulation des lettres i/j et u/v.

Les autres ouvrages appartiennent à la catégorie générique « typographie moderne ».



## EXEMPLE



<http://gallica.bnf.fr/ark:/12148/bpt6k6547544k>

Autres langues romanes et langues écrites avec les alphabets grec et latin

Comme il a été dit section 4.1.2, ces contenus sont à traiter en OCR brut.

Mélange de langues et/ou d'alphabets

Seuls les contenus en français présents dans des blocs mixtes (mélangeant langues et/ou alphabets) seront pris en compte dans l'évaluation du taux OCR qualité garantie.

*Taux OCR qualité éditoriale*

Selon les prestations, ce taux pourra être exigé sur des zones spécifiques des documents (par exemples les titres). Il diffère des taux qualité garantie ou haute qualité par les règles suivantes :

- périmètre :
  - tous les blocs de texte en alphabet latin,
  - seuls les mots ou blocs illisibles sont exclus,

- jeu de caractères :
  - tous les caractères ASCII font partie du jeu de caractères garanti
  - la casse et les caractères accentués sont pris en compte.

#### 4.1.4 Encodage

La transcription du texte dans les fichiers ALTO se fera avec l'encodage Unicode UTF-8 restreint à l'ensemble des blocs de caractères Unicode nécessaires à la transcription des langues romanes plus les blocs nécessaires à la transcription des caractères grecs.

#### 4.1.5 Signes typographiques, caractères et symboles spéciaux, etc.

##### *Signes typographiques*

Les tirets d'incise, d'énumération et de liste (en général un tiret demi-cadratin) sont à produire dans le flot de texte avec le caractère Unicode &#x2013;

Les tirets de dialogue, les tirets longs (en général un tiret cadratin) sont à produire dans le flot de texte avec le caractère Unicode &#x2014;

Les puces de forme standard (•) sont à produire dans le flot de texte avec le code &#x2022;

##### *Caractères spéciaux*

Les caractères spéciaux sont traités en OCR et intégrés dans le flux texte sans traitement particulier. Ils peuvent entraîner la génération de mots illisibles.



#### **EXEMPLE**

BÉLÈGUÉ

Président : M. Lucien CORNET, 10, rue de l'Écrivain, à Sens.  
 Vice Présidents { M. Désiré BUDAN ☞, propriétaire, à Wo-l'Archevêque.  
                           M. Georges RAVIN ☞, propriétaire, 8, rue des  
                           Francs-Bourgeois, à Sens.  
 Secrétaire : M. Désiré GORCE ☞, propriétaire, à Courtois.  
 Trésorier : M. COCHARD Henri, propriétaire, à Collemiers.  
 Secrétaire-adjoint : M. AUPIERRE ☞, propriétaire à Maillot.  
 Trésorier-adjoint : M. PRIMAULT Anatole ☞, propriét., à St-Clément.  
 Bibliothécaire-archiviste : M. ROGNON Désiré ☞, prop., à Chaumont.

**BÉLÈGUÉ DE PARIS**



#### **NOTES**

Comme mentionné section 4.1.3, ces caractères ne sont pas comptabilisés dans l'évaluation du taux OCR qualité garantie.

### *Points de suite*

Les points de suite, notamment dans les tables des matières et les index, peuvent être absents du flux de texte produit.

## 4.2 Styles typographiques

Par page, la liste des styles de paragraphes et des styles de caractères utilisés est donnée dans l'en-tête du fichier ALTO (élément <Styles>).

Au niveau de chaque bloc de texte, seul le style majoritaire est indiqué.

Au niveau de chaque ligne, seul le style majoritaire est indiqué.

Au niveau de chaque mot, l'éventuel style utilisé est indiqué. Un seul style est possible par mot.

### 4.2.1 Niveau de qualité

Les styles sont obtenus par un traitement purement automatique lors de l'OCR brut.



Alors que les exigences de validité et de synchronisation des balises de styles sont respectées scrupuleusement (taux qualité garantie), l'exactitude des styles (police, taille, enrichissements) en regard du document d'origine ne peut pas l'être. Ainsi, la reconnaissance des styles n'est pas soumise au taux qualité garantie.

### 4.2.2 Exposant, indice

Les mots composés en exposant ou indice seront identifiés par le mécanisme des styles.

### 4.2.3 Titres

Les titres ne sont pas identifiés en tant que tel, sauf si une tâche de reconnaissance des titres est demandée (cf. section 5.6.2).

### 4.2.4 Typographies mal reconnues par l'OCR

Différentes typographies et écritures sont mal reconnues par l'OCR :

- écriture manuscrite,
- police fantaisie, script, avec relief, double traits, etc.
- police italique avec un fort degré d'inclinaison (supérieur à 30°),
- police de type Fraktur, gothique, qui requiert des licences OCR spécifiques,
- alphabet non latin (russe, asiatique, etc.).



## NOTES

Ces portions de texte mal reconnues par l'OCR seront considérées comme étant « illisibles » (*illegibles*) et seront traitées en OCR brut. Ce processus est détaillé aux sections 5.6.9 et 6.2.2.

### *Ecriture manuscrite*

Ce cas est décrit aux sections 5.7.5 et 5.8.5.

### *Polices script*

Les blocs de texte composés avec une police script sont décrits avec un élément TextBlock doté d'une étiquette LayoutTag LABEL="scriptFonts" (cf. section 5.10).



## EXEMPLE

*Platinum  
Gifts for Her  
Precious Gems  
Legacy Collection  
Etoile Diamond & Silver*

---

### *Autre cas de polices non exploitables*

Ces blocs de texte sont décrits avec un élément TextBlock doté d'une étiquette LayoutTag LABEL="illegible" (cf. section 5.10).

 **EXEMPLE**

| Patable.   |  |
|--|--|
| <b>Comence la table de ce present liure.</b>                                 | De rondellis de exercis specie et alia de quatuordecim quinq; silabarum et de quatuor sillabarum et alia multa que sequuntur per octid. m. fucillet. viij. |
| ¶ Et puenfremet le prologue  | De plures specie et calibus/sez quater. v. vi. et de septimo colore de forma ptale. modic. fucillet. iij.  |
| Diffinitio puma et leuandit capitulum fucillet.                              | De octo: a: n: o: n: a: de c: i: n: a: d: e: c: i: a: d: u: o: s: decena specie. fucillet. x. et. vij.   |
| De dinatione rethorice. fucillet.  | De terdecimo colore rethorice de sermas casis de. xiiij. xv. xvi. xvij. specie dno et  |
| De specificatione quad: in uet: et pua: la de rethorica. iij. cap. fucillet. | De iij. per ordinatis. fucillet. x. vij. et. xiiij.  |
| De uetis. iij. capitula. fucillet.   | Decimum caplu persona compilandi mcaulatis. fucillet. xj.  |
| De uetis uicoune. fucillet.  | De non edis. fucillet. iij.  |
| De figuris. v. capitula. fucillet.   | De milleris compilatis cronis uenit os et hiltouis. fucillet. xij.   |
| De definitione figurarum. fucillet.  | ¶ Preface. fucillet. xiiij.  |
| Diffinitio finalis pbe. p:is notabile cau: la uilitatis. fucillet.           | ¶ La doctore de megere. fucillet. xiiij.   |
| Diffinitio saxope secundu: aliud notabile. fucillet.                         | ¶ De lorde royal. fucillet. xiiij.   |
| De apocopa textis. fucillet.   | ¶ De salet michel archage. fucillet. xiiij.  |
| ¶ Ansa uilitatis apocopa. fucillet.  | ¶ De donner baillie sur feu roy charles. huy: leuc de ce uenit. fucillet. xj.  |
| De linoninis quarto. fucillet.   |  |
| ¶ De diffinitione equinoctialis.   |  |

*Alphabets non latins*

Ce cas est décrit à la section 4.1.2.

**4.3 Césures**

La partie reconnue est dans l'attribut content, complétée par les attributs :

- subs\_type : précise la partie concernée, HypPart1 pour la première partie du mot césuré, HypPart2 pour la seconde ;
- subs\_content : restitue le mot complet non césuré.

L'élément HYP représente le tiret de césure.

 **EXEMPLE**

```

...
<String ID="PAG_00000001_ST000100" STYLEREFs="TXT_1" HPOS="1285"
VPOS="1910" HEIGHT="47" WIDTH="76" WC="0.99" CONTENT="va-"
SUBS_TYPE="HypPart1" SUBS_CONTENT="variées,"/>
<HYP HPOS="1361" VPOS="1957" WIDTH="30" CONTENT="-"/>
</TextLine>

<TextLine ID="PAG_00000001_TL000028" STYLEREFs="TXT_1" HPOS="336"
VPOS="1955" HEIGHT="66" WIDTH="1024">
<String ID="PAG_00000001_ST000101" STYLEREFs="TXT_1" HPOS="336"
VPOS="1984" HEIGHT="37" WIDTH="127" WC="0.99" CONTENT="riées,"
SUBS_TYPE="HypPart2" SUBS_CONTENT="variées,"/>

```



1. Dans le cas d'une double césure (par ex. « a-bracada-bra »), on utilisera un SUBS\_TYPE="HypPart1" suivi de deux SUBS\_TYPE="HypPart2".

2. Il n'est pas demandé de distinguer entre tiret de césure et trait d'union. Un mot composé placé sur deux lignes pourra donc être traité avec une césure.

#### 4.3.1 Répartition entre TextBlock

Le prestataire veillera à éviter de répartir un mot césuré sur deux TextBlock (sauf dans le cas où la césure intervient sur un mot à cheval sur deux colonnes ou deux pages).

## 5. SEGMENTATION ET STRUCTURATION

---

Ce chapitre décrit de quelle manière la structure du document sera représentée dans le fichier ALTO produit en sortie.

Un taux qualité spécifique au marché (cf. CCTP) s'applique à cette segmentation, indépendamment du taux qualité OCR (cf. section 6).



SAUF MENTION CONTRAIRE, TOUS LES ELEMENTS DECRITS DANS CE CHAPITRE SONT SOUMIS AU TAUX QUALITE GARANTIE.

### 5.1 Description des pages

#### 5.1.1 Numérotation des pages

L'attribut PHYSICAL\_IMG\_NR (obligatoire) de l'élément Page permet de renseigner le numéro de la page dans l'ordre séquentiel des pages du document.

L'attribut ID de l'élément Page représente le numéro d'ordre dans la numérotation des fichiers ALTO :

```
<Layout>
  <Page ID="PAG_00000002" HEIGHT="3353" WIDTH="2065"
    PHYSICAL_IMG_NR="2" > ...
</Page>
</Layout>
```



CET ATTRIBUT EST OBLIGATOIRE.

L'attribut PRINTED\_IMG\_NR (optionnel) de l'élément Page permet de renseigner les informations de page/folio ou les éléments descriptifs de couverture/couverture, tel qu'imprimé sur la page (cf. « Référentiel d'enrichissement des métadonnées », section 3.3 pour les règles de valorisation de cet attribut).

```
<Layout>
  <Page ID="PAG_00000135" HEIGHT="3353" WIDTH="2065"
    PHYSICAL_IMG_NR="135" PRINTED_IMG_NR="129" > ...
</Page>
</Layout>
```



Quand le fichier ALTO représente une double page, on utilisera la syntaxe x-y pour ces deux attributs. Pour une double page de folio 8 et 9, on aura donc :

```
<Layout>
```

```
<Page ID="PAG_000008" HEIGHT="3353" WIDTH="4050"
  PHYSICAL_IMG_NR="8-9" PRINTED_IMG_NR="8-9"> ...
</Page>
</Layout>
```

### 5.1.2 Qualité de numérisation des pages

L'attribut QUALITY (obligatoire) de l'élément Page permet d'indiquer au lecteur l'état du document d'origine et en conséquence son résultat numérique. Cette information correspond aux commentaires de type « usager » portés au niveau du feuillet ou de la page concernée dans l'élément <vueObjet> du refNum (cf. « Référentiel d'enrichissement des métadonnées », section 4.1).

Sa valeur sera "OK" dans le cas courant (rien à signaler) ou "Damaged" : tous les autres cas tel que décrit dans le « Référentiel d'enrichissement des métadonnées ».

Dans ce dernier cas, l'attribut QUALITY\_DETAIL de l'élément Page permet de préciser le type de défaut. Il contiendra la même valeur que l'élément <commentaire> de l'élément <vueObjet> du refNum :

**ALTO :**

```
<Layout>
  <Page ID="PAG_00000101" HEIGHT="3353" WIDTH="2065"
    PHYSICAL_IMG_NR="101" PRINTED_IMG_NR="94"
    QUALITY="damaged" QUALITY_DETAIL="Report d'encre"> ...
  </Page>
</Layout>
```

**refNum :**

```
<commentaire type="USAGER" date="2013-09-30T15:00:43Z">Report d'encre
</commentaire>
```

### 5.1.3 Qualité de l'océrisation des pages

L'attribut ACCURACY de l'élément Page permet d'indiquer le taux OCR estimé pour la page, exprimé en %. Cette valeur doit être comprise entre 0.0 et 100.0.

```
<Layout>
  <Page ID="PAG_00000135" HEIGHT="3353" WIDTH="2065"
    PHYSICAL_IMG_NR="135" PRINTED_IMG_NR="129"
    ACCURACY="98.380"> ...
```

### 5.1.4 Pages vides

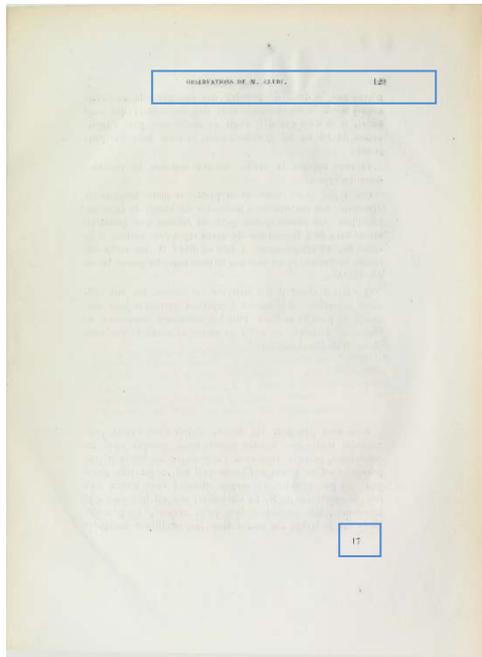
L'attribut PAGECLASS de l'élément Page permet de décrire certains cas particuliers.

Les pages vides (sans aucun contenu) sont décrites avec PAGECLASS="BlankPage".

```
<Layout>
  <Page ID="PAG_00000002" HEIGHT="3353" WIDTH="2065"
    PHYSICAL_IMG_NR="1" QUALITY="OK" PAGECLASS="BlankPage">
  </Page> ...
</Layout>
```

### 5.1.5 Pages avec contenu en marges mais sans contenu principal

Si une page présente un contenu à placer en XxxMargin, mais que le reste de la page est vide, il faut utiliser la description PAGECLASS="BlankPrintSpace".



```
<Layout>
  <Page ID="PAG_00000135" PAGECLASS="BlankPrintSpace" HEIGHT="4153"
    WIDTH="3049" PHYSICAL_IMG_NR="135" PRINTED_IMG_NR="129"
    QUALITY="OK" > ...

  <TopMargin ID="PAG_00000135_TopMargin" HPOS="0" VPOS="0" HEIGHT="562"
  WIDTH="3049">
    <TextBlock ID="PAG_00000135_TB000001" STYLEREFs="TXT_1" HPOS="1124"
    VPOS="517" HEIGHT="44" WIDTH="1224" LANG="fr">
      <TextLine ... </TextLine>
    </TextBlock>
  </TopMargin>

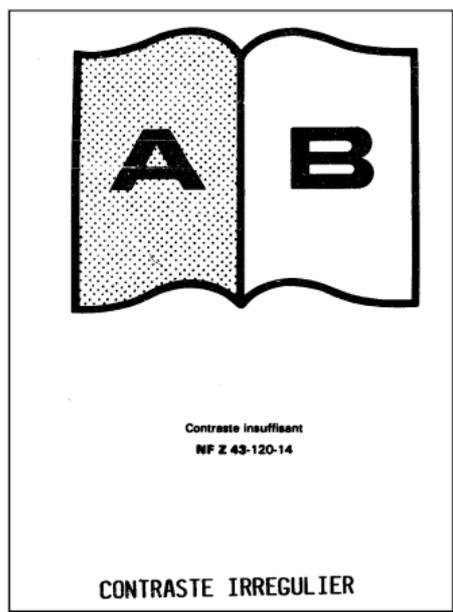
  <BottomMargin ID="PAG_00000135_BottomMargin" HPOS="0" VPOS="3363"
  HEIGHT="790" WIDTH="3049">
    <TextBlock ID="PAG_00000135_TB000002" STYLEREFs="TXT_1" HPOS="2142"
    VPOS="3363" HEIGHT="40" WIDTH="46" LANG="fr">
      <TextLine ... </TextLine>
    </TextBlock>
  </BottomMargin>
</Page>
</Layout>
```

---

### 5.1.6 Pages de logo

Les pages ne contenant qu'un logo (page de type « L » dans le manifeste numérique refNum) sont traitées comme des pages blanches.

Les éventuelles légendes textuelles en plus de la représentation graphique à proprement parler ne seront pas traitées.



Logo + légende + n° NF du logo + légende

## 5.2 Orientation de la page

Les images numérisées seront toujours fournies dans le sens de l'original.

L'orientation de la lecture se détermine sur la totalité de la page et non pas seulement sur la partie textuelle.

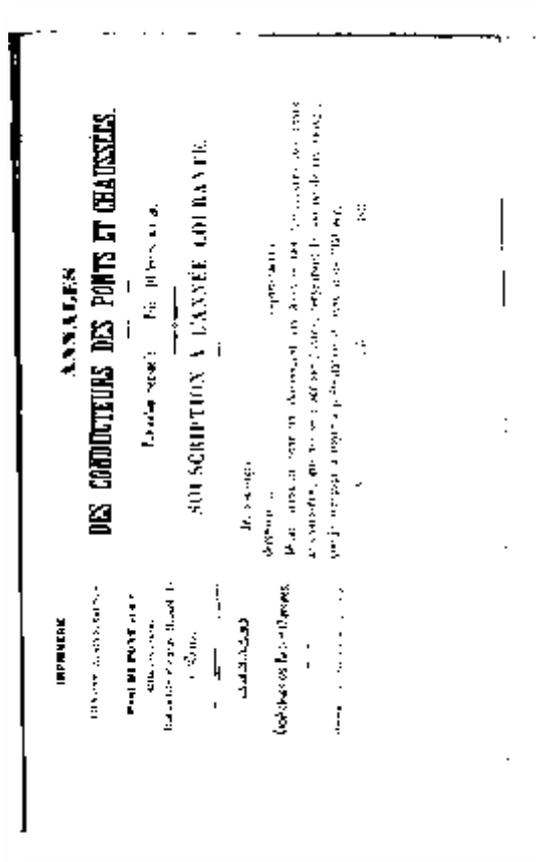
Si l'orientation de l'original (c'est-à-dire celle de la reliure) est différente de l'orientation de la lecture, l'angle de rotation des éléments de la page doit être spécifié dans le fichier ALTO à l'aide de l'attribut Rotation, selon un angle exprimé en degrés, dans le sens inverse des aiguilles d'une montre (sens trigonométrique).

Deux cas peuvent se présenter :

- Si tous les éléments de la page sont dans une orientation de lecture donnée, alors en restitution, ces éléments comporteront un attribut Rotation. Dans l'exemple suivant, les composants de la page sont produits avec un attribut rotation à 90° :



```
<TextBlock ID="PAG_248_TB000001" STYLEREFS="TXT_1" HPOS="1996"  
VPOS="1019" HEIGHT="1557" WIDTH="66" LANG="fr" ROTATION="90">
```



*Sens de la page originale : 90°*

- Si la page comprend des zones ayant plusieurs orientations de lecture, on détermine l'orientation de lecture majoritaire selon l'ordre de priorité suivant :
  - L'entête et le numéro de page ne sont pas pris en compte (ils sont structurés en TopMargin).
  - Les éléments GraphicalElement ne sont pas pris en compte.
  - Pour déterminer le sens de lecture, le « corps du texte » prime sur le reste (les illustrations, les tableaux avec leurs éventuels contenus textuels et les éléments graphiques, etc.) s'il occupe plus d'un tiers de la page.
  - De la même manière, le « corps du texte » prime sur les légendes des illustrations et des tableaux, même quand ceux-ci sont en dehors de l'objet Illustration.

Une fois l'orientation de lecture majoritaire déterminée, on applique l'attribut Rotation aux zones dont l'orientation de lecture est différente de l'orientation de lecture majoritaire.



## 5.3 Structuration de la page

### 5.3.1 PrintSpace

Le contenu spécifique à la page doit être inclus dans le PrintSpace. En cas de doute, il est possible d'inclure tout le texte d'une page dans le PrintSpace.

Le PrintSpace peut être serré autour du texte (c'est-à-dire ne pas englober l'ensemble de la page) à condition que tout le texte de la page (hors marge) soit compris dans celui-ci.



#### **ATTENTION**

LES PAGES BLANCHES DOIVENT ETRE DECRITES SELON LES REGLES EXPOSEES A LA SECTION 5.1.1.

### 5.3.2 XxxMargin

Le contenu récurrent et répétitif (titre, intitulé d'une section, sous-titre, nom d'auteur et numéro de page) doit appartenir à un xxxMargin.

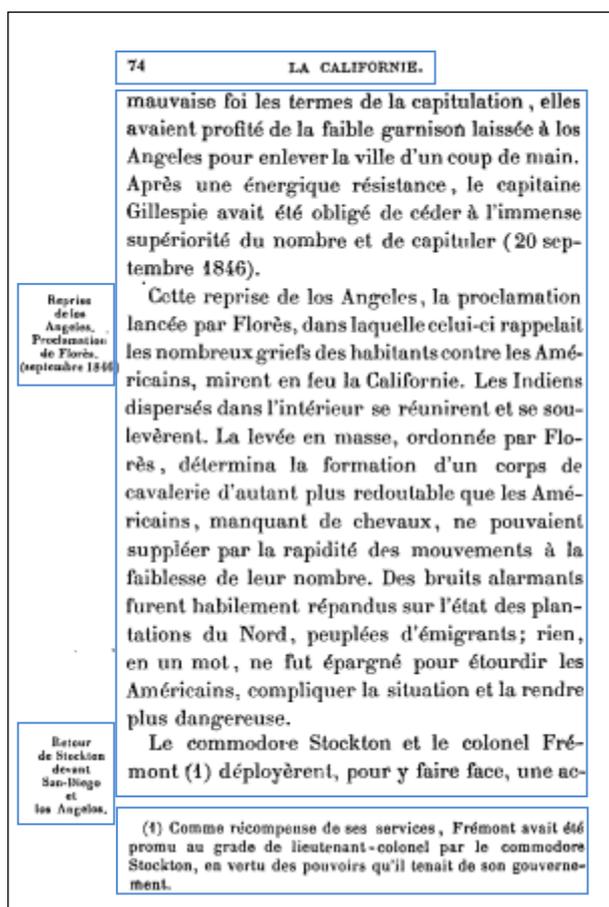
Si un XxxMargin est indiqué, il peut couvrir tout l'espace entre le bord de la page concernée et le(s) printSpace concerné(s).

Un contenu non répétitif et spécifique à la page (notes ou en-têtes de section en marge par exemple) fait partie du printSpace.



#### **ATTENTION**

UNE MISE EN PAGE AVEC DES BLOCS EN MARGE A DROITE OU EN MARGE A GAUCHE NE SIGNIFIE PAS NECESSAIREMENT QUE CES BLOCS DOIVENT ETRE PLACES EN LEFTMARGIN OU EN RIGHTMARGIN.



Le titre courant (« 74 LA CALIFORNIE. ») doit constituer un topMargin.

PrintSpace : le contenu des deux notes marginales à gauche est lié au corps du texte de la page et il lui est spécifique. Il en ressort qu'ils ne font pas partie d'un leftMargin mais doivent être séparés afin que des parties de leur contenu ne soient pas mélangées avec le texte de la colonne principale.

Dans cet exemple, les notes marginales apparaitront dans le flux ALTO après les notes de bas de page.

## 5.4 Blocs manqués

Les blocs de texte non identifiés par la segmentation automatique seront identifiés visuellement et décrits manuellement : position, taille, type.

Les blocs texte ainsi identifiés seront ensuite traités en OCR. Si le moteur OCR ne détecte aucun texte à l'intérieur d'un tel bloc, le bloc texte sera rempli avec une seule ligne composée de la mention "[texte manquant]", et il sera typé avec une étiquette LayoutTag LABEL="missing" (cf. section 5.10).



L'identification visuelle des blocs manqués n'est pas soumise au taux qualité garantie.

## 5.5 Ordre de lecture et ordre des segments

L'ordre de lecture est également déterminé lors de la reprise de la segmentation automatique, avec application des règles qui suivent :

- Le texte se lit en colonne, de haut en bas et de gauche à droite.
- Les illustrations sont segmentées dans l'ordre où elles apparaissent, de haut en bas et de gauche à droite.
- Quand une illustration s'accompagne d'une légende, utiliser un ComposedBlock qui inclut l'image et la zone de texte.

Dans les fichiers XML ALTO, l'ordre de lecture est indiqué par la numérotation séquentielle des blocs, via leur attribut ID, ainsi que par leur ordre de présence dans le fichier XML.

Cette numérotation séquentielle est propre à chaque type de bloc. Par exemple, les éléments String auront une numérotation de PAG\_00000001\_SP000001, PAG\_00000001\_SP000002 à PAG\_00000001\_SP00000n, les éléments TextBlock de PAG\_00000001\_TB000001, PAG\_00000001\_TB000002 à PAG\_00000001\_TB00000n, etc.

Dans chaque identifiant, le premier numéro est le numéro de page, et le second l'ordre de l'élément dans sa séquence.



L'ordre de lecture au sein des zones de marges suit les mêmes règles.

### 5.5.1 Mise en page en colonnes

En règle générale :

- Les entêtes des colonnes, si elles ne sont pas dans le topMargin, figurent en haut de la colonne la plus proche.
- Les bas des colonnes, si elles ne sont pas dans le bottomMargin, figurent en bas de la colonne la plus proche.
- Les colonnes sont traitées dans le sens de lecture.
- La totalité des segments de la colonne N sont traités avant un segment qui appartient à une colonne N+1.
- Les filets verticaux de gouttière ne sont pas segmentés.



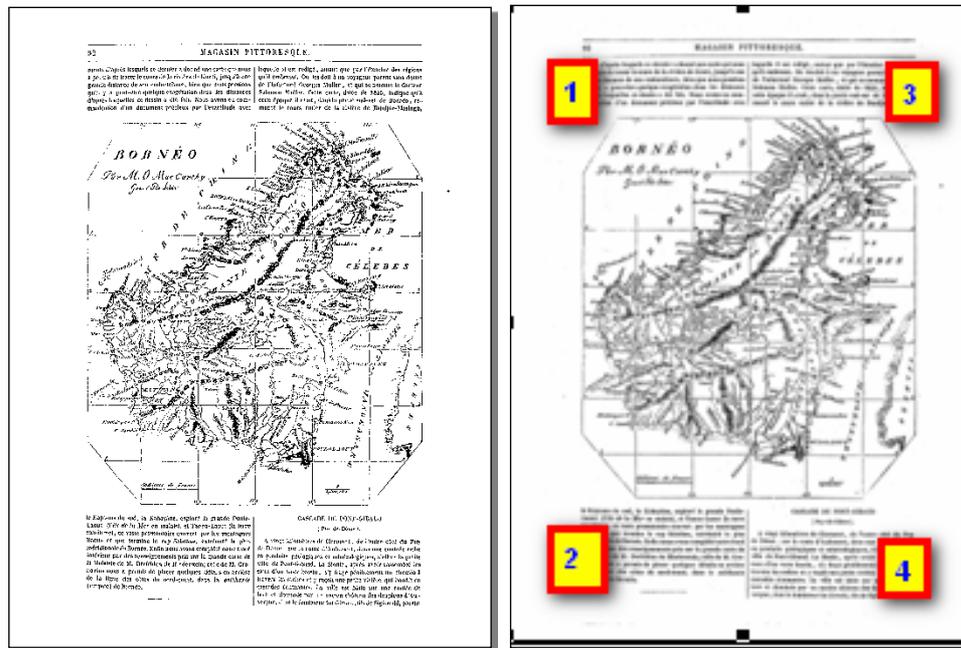
Double colonne type imprimé

## 5.5.2 Mise en page en colonnes avec des éléments centraux non textuels

Un élément central non textuel est un élément non textuel (illustration, carte etc.) qui occupe (pleinement ou partiellement) plusieurs colonnes.

En règle générale :

- L'ordre de lecture suit les règles générales sur les colonnes.
- Le bloc central non textuel peut être mis à n'importe que point entre les TextBlock ainsi ordonnés.



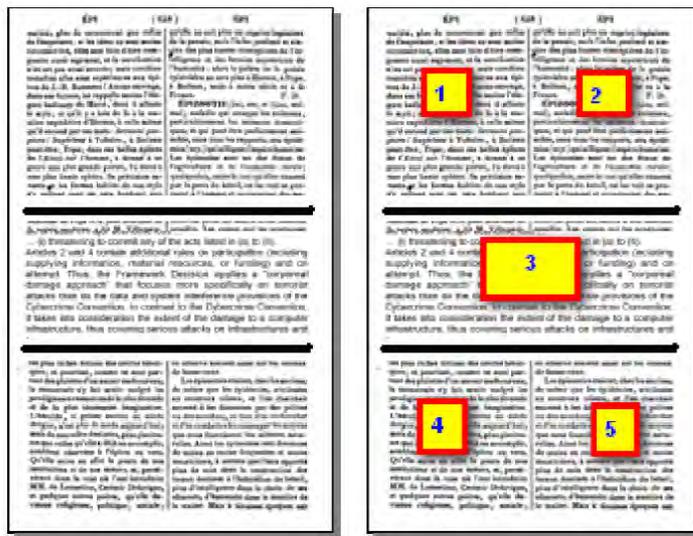
*La légende de l'illustration n'est pas centralisée, elle sera incluse dans « sa » colonne*

### 5.5.3 Mise en page en colonnes avec des éléments centraux textuels

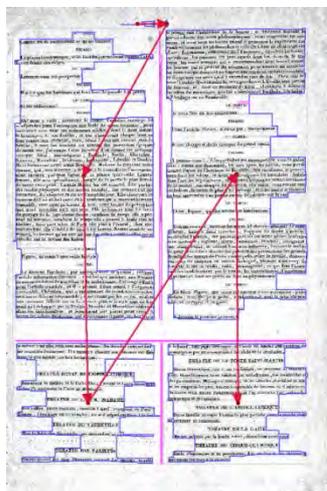
Un élément central textuel est un élément textuel (simple texte, encadré avec texte, table, etc.) qui occupe (pleinement ou partiellement) plusieurs colonnes.

En règle générale :

- En plus des règles sur les colonnes, les blocs dont la position verticale est avant un bloc textuel central, figurent avant ce bloc textuel central.



Ce cas est particulièrement courant pour les contenus presse et peut conduire à un ordre de lecture erroné.



Ordre de lecture erroné

## Mise en page avec imbrication de colonnes

Les mises en page de type presse présentent des particularités :

- La zone de l'ours sera toujours placée en premier dans l'ordre de lecture, même si elle ne court pas sur toute la largeur de la page.



- L'ordre de lecture sera déterminé en appliquant les règles générales d'abord sur les zones de plus haut niveau, puis à l'intérieur de chaque zone.



Identification des zones de plus haut niveau

<http://gallica.bnf.fr/ark:/12148/bpt6k502978s>



**Une année...**

**L'Or**

**Un message de Roosevelt au Congrès**

**LES NATIONS ALIMENTAIRES DE JANVIER**

**Le Maréchal a reçu les enfants de France**

**Les mystères du marché noir**

**Un wagon de dix tonnes de cacahuètes était égaré à l'année (Mayenne)**

**Maux tabaciques qui cherchent à diverger sur l'air**

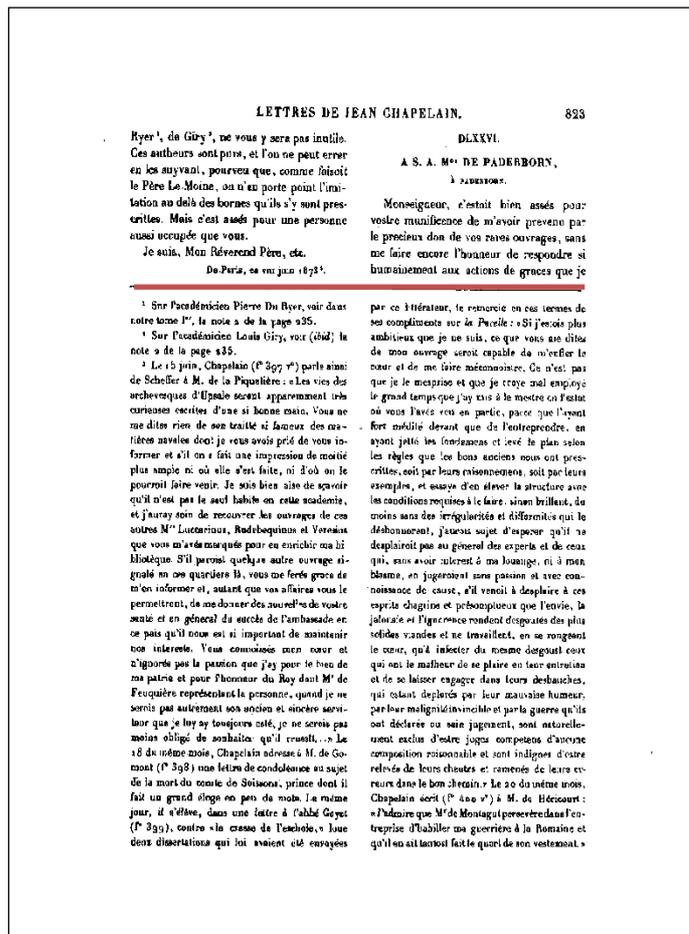
**Les troupes anglo-gaullistes auraient occupé Djibouti**

Ordre de lecture interne dans la zone n° 2

## 5.5.4 Corps du texte et notes séparés par un trait

La séparation entre le corps du texte et les notes sera contrôlée.

Si aucun trait séparateur ou indice de séparation n'est apparent, alors le traitement sera basé sur l'exploitation de deux colonnes « simples ».



Exemple de notes de bas de page composées sur deux colonnes

## 5.6 Texte

### 5.6.1 Paragraphes

Lors de l'étape de segmentation des ouvrages, les blocs textes faisant partie du PrintSpace sont détectés sans séparation entre les paragraphes. Ils sont ensuite envoyés dans l'étape d'OCR.

Si les blocs textuels sont séparés par des éléments graphiques ou des tableaux, plusieurs blocs de texte sont segmentés.

Après l'OCR, c'est lors de l'étape de structuration (processus automatique avec validation manuelle) que le découpage du corps de texte en paragraphes a lieu et que les paragraphes sont identifiés.



La BnF attire l'attention du prestataire sur la nécessité d'un découpage en paragraphes le plus fidèle possible à l'original, afin de permettre un bon reformatage des textes lors de leur diffusion sur le Web ou sur supports nomades.

## 5.6.2 Titres

### *Titres*

Selon la nature des contenus à numériser, la prestation pourra inclure une reconnaissance des titres, notamment dans le cas de documents de type presse (titres et rubriques d'article de journaux).

Dans le cas où une détection de titres est demandée, le niveau de titre est stocké par une étiquette StructureTag sur le TextBlock ou le TextLine concerné (cf. section 5.6.2 et 5.10).

La tâche de reconnaissance des titres est décrite en détail dans le « Référentiel d'enrichissement du texte ».

### *Titres de presse*

Les titres de presse ainsi que les ours pourront être déqualifiés en OCR brut, du fait des polices ou des tailles de caractères utilisées.



## 5.6.3 Tableaux

Les tableaux sont identifiés lors de l'étape de segmentation.



Sont définis comme tableaux les éléments multicolonnés qui font une rupture dans le flux de texte, ayant ou non un cadre externe et des traits de séparation entre les cellules.

Les tableaux contenant des chaînes alphanumériques seront convertis en texte non structuré en tableau, en utilisant une étiquette LayoutTag LABEL="table" sur

l'élément TextBlock (cf. section 5.10). L'ordre de lecture à l'intérieur des différentes parties des tableaux est celui obtenu par le moteur OCR, c'est-à-dire en lignes d'abord, puis en colonnes.

Chaque bloc de texte contenu dans le tableau sera marqué comme un TextBlock avec une étiquette LayoutTag LABEL="table". Il peut y avoir plusieurs blocs de texte de type Table pour un même tableau physique sur une même page.



**NOTES**

Les éléments graphiques au sein d'un tableau qui font office de filets ou de traits de séparation ne seront pas segmentés.

| DESIGNATION          | MOYENS de piles | MOYENS de piles de piles | MOYENS de piles de piles | DIAMÈTRE de piles | PROFONDEUR de piles | SUMMAIRES payés aux ouvriers | NATURE DES COUPURES à traverser  | OBSERVATIONS  |
|----------------------|-----------------|--------------------------|--------------------------|-------------------|---------------------|---------------------|---------------------|---------------------|---------------------|------------------------------|--|---|
| Pile de Strood...    | 14              | 5                        | 2.456                    | 18.94             | 13.20               | 140.00              | 5,319               | 34                  |                     |                              | Débris de maçonnerie et de charpente, débris de maçonnerie et de charpente.<br>Terrain naturel, c'est-à-dire sable, gravier et cailloux.<br>Pierre, débris de charpente et craie, terrain naturel. | Entoncé hivernal, sans air comprimé.<br>Entoncé hivernal sans air comprimé. |
| Pile de Rochester... | 14              | 5                        | 2.134                    | 15.42             | 6.71                | 89.60               | 6,587               | 50                  |                     |                              |  |   |
| Culée de Strood...   | 30              | 2                        | 1.45                     | *                 | 5.49                | 84.30               | 4,307               | 50                  |                     |                              |  |   |
| Culée de Rochester   | 6               | 2                        | 1.45                     | *                 | 5.49                | 40.60               | 1,035               | 54                  |                     |                              |  |   |
|                      | 3               | 2                        | 1.85                     | 13.81             | 6.25                | 70.58               | 2,620               |                     |                     |                              |  |   |

A (1) L'unité est le jour de 10 heures.  
 B (2) Ces sommes sont seulement celles payées pour la main d'œuvre d'entretien et de détail, non compris les frais pour travaux accessoires, échafaudage, etc.

Eléments non segmentés

Tableau avec filet de séparation et accolades

et le Canada en étaient les principaux destinataires; depuis 1935, au contraire, c'est le groupe Belgique-Allemagne-Pays-Bas qui figure en tête de la liste des importateurs :

destinations déclarées des cargaisons « uniques » de manganèse

| année                         | Angleterre | Etats-Unis Canada | Belgique Allemagne Pays-Bas | France | autres |
|-------------------------------|------------|-------------------|-----------------------------|--------|--------|
| (en milliers de tonnes poids) |            |                   |                             |        |        |
| 1911..                        | 22         | 55                | 10                          | —      | 7      |
| 1912..                        | 80         | 84                | 39                          | 19     | 12     |
| 1913..                        | 106        | 85                | 64                          | 13     | 54     |
| 1920..                        | 33         | 18                | 37                          | 21     | 27     |
| 1921..                        | 22         | 30                | 61                          | 28     | 4      |
| 1922..                        | 20         | 5                 | 136                         | 21     | 12     |
| 1923..                        | 44         | 6                 | 74                          | 20     | 9      |
| 1924..                        | 34         | —                 | 92                          | 24     | 2      |
| 1925..                        | 36         | 25                | 58                          | 34     | 13     |
| 1926..                        | 5          | 12                | 28                          | 7      | 16     |
| 1927..                        | 30         | 7                 | 122                         | 35     | 28     |
| 1928..                        | 31         | 12                | 82                          | 28     | 12     |
| 1929..                        | 105        | 74                | 105                         | 37     | 25     |
| 1930..                        | 100        | 10                | 64                          | 41     | 24     |
| 1931..                        | 17         | —                 | 38                          | 25     | 40     |
| 1932..                        | 20         | —                 | 16                          | 10     | 6      |
| 1933..                        | 12         | —                 | 16                          | 14     | 7      |
| 1934..                        | 47         | —                 | 25                          | 11     | 20     |
| 1935..                        | 45         | —                 | 104                         | —      | 25     |
| 1936..                        | 20         | 54                | 74                          | —      | 22     |
| 1937..                        | 128        | 48                | 217                         | 104    | 72     |

En bref, le groupe Belgique-Allemagne-Pays-Bas, qui avait importé 10,6 % du total des cargaisons uniques en 1910, puis 18,2 % en 1913, s'inscrivit, l'année dernière, pour 36,5 %.

Pour les produits autres que le manganèse, des modifications sensibles sont également à noter, en particulier au cours de la période 1929-1937, ainsi qu'il ressort du tableau suivant :

principaux produits miniers, autres que le manganèse

| année                         | zinc | plomb | fer | cuivre | houille |
|-------------------------------|------|-------|-----|--------|---------|
| (en milliers de tonnes poids) |      |       |     |        |         |
| 1926..                        | 320  | 216   | 111 | 19     | —       |
| 1927..                        | 204  | 223   | 222 | 28     | —       |
| 1928..                        | 158  | 307   | 137 | 22     | —       |
| 1929..                        | 106  | 200   | 125 | 35     | —       |
| 1930..                        | 127  | 213   | 161 | 53     | —       |
| 1931..                        | 127  | 216   | 155 | 61     | 2       |
| 1932..                        | 164  | 126   | 116 | 47     | 1       |
| 1933..                        | 90   | 90    | 206 | 40     | 105     |
| 1934..                        | 93   | 82    | 205 | 106    | 105     |

La déclin, très parallèle, des deux trafics du zinc et du plomb en 1936 et 1937 est en grande partie imputable au fait que, depuis deux ans, une fraction importante de ces exportations originaires d'Australie échappe momentanément au Canal :

envois d'Australie

| année                      | zinc | plomb |
|----------------------------|------|-------|
| (milliers de tonnes poids) |      |       |
| 1929..                     | 94   | 125   |
| 1930..                     | 86   | 129   |
| 1931..                     | 59   | 129   |
| 1932..                     | 10   | 149   |
| 1933..                     | 43   | 120   |
| 1934..                     | 47   | 179   |
| 1935..                     | 51   | 154   |
| 1936..                     | 4    | 24    |
| 1937..                     | 19   | 21    |

Fort heureusement des compensations, une fois de plus, ont joué; et les pertes constatées aux rubriques du zinc et du plomb se sont trouvées contre-balançées, et au delà, par l'importance sans précédent qu'ont connue en 1937 les envois de soufre indienne, les expéditions d'aminite de même provenance, ainsi que les passages de houille embarquée aux Indes Néerlandaises.

Rappel des recettes quotidiennes de Avril 1937 pour comparaison avec Avril 1936

| Livres Sterling               |        | Livres Sterling               |        | Livres Sterling               |        |
|-------------------------------|--------|-------------------------------|--------|-------------------------------|--------|
| rapet. 1936                   |        | rapet. 1937                   |        | rapet. 1937                   |        |
| 1 <sup>er</sup> .....         | 71.000 | 14.....                       | 65.400 | 21.....                       | 64.700 |
| 2.....                        | 71.000 | 15.....                       | 65.300 | 22.....                       | 64.700 |
| 3.....                        | 71.000 | 16.....                       | 65.300 | 23.....                       | 64.700 |
| 4.....                        | 71.000 | 17.....                       | 65.300 | 24.....                       | 64.700 |
| 5.....                        | 71.000 | 18.....                       | 65.300 | 25.....                       | 64.700 |
| 6.....                        | 71.000 | 19.....                       | 65.300 | 26.....                       | 64.700 |
| 7.....                        | 71.000 | 20.....                       | 65.300 | 27.....                       | 64.700 |
| 8.....                        | 71.000 | 21.....                       | 65.300 | 28.....                       | 64.700 |
| 9.....                        | 71.000 | 22.....                       | 65.300 | 29.....                       | 64.700 |
| 10.....                       | 71.000 | 23.....                       | 65.300 | 30.....                       | 64.700 |
| Du 1 <sup>er</sup> au 10..... |        | Du 1 <sup>er</sup> au 10..... |        | Du 1 <sup>er</sup> au 30..... |        |
| 283.000                       |        | 283.000                       |        | 283.000                       |        |

OPPOSITIONS

Obligations 5 % 3<sup>e</sup> série

Sur le 10<sup>e</sup> coupon, à l'échéance du 1<sup>er</sup> mars 1938.

N<sup>o</sup> 85.022 — 122.628 — 206.029 — 263.630

270.229 (opposition n<sup>o</sup> 9.816).

Parts de Fondation

N<sup>o</sup> 9.259 (opposition n<sup>o</sup> 3.817).

Tableaux sans filet de séparation, dans une maquette en deux colonnes

### Taux qualité des tableaux

Le taux qualité des contenus tableau sera traité différemment pour les deux types de tableau suivants :

- tableau simple et court de deux colonnes, le contenu étant généralement :
  - composé avec le même corps que le texte courant,
  - présenté sans filet ni bordure,
  - les deux colonnes étant séparées par des tabulations ou des points de suite.

Ces contenus seront traités en OCR **taux qualité garantie**. Ils seront balisés :

**TAGREFS="TAG\_table"**

**ABRÉVIATIONS**

|                           |                                    |                        |                        |
|---------------------------|------------------------------------|------------------------|------------------------|
| <i>alt.</i> . . . . .     | altitude.                          | <i>long.</i> . . . . . | longueur.              |
| <i>aub.</i> . . . . .     | auberge.                           | <i>Lun.</i> . . . . .  | Lundi.                 |
| <i>auj.</i> . . . . .     | aujourd'hui.                       | <i>m.</i> . . . . .    | mètre.                 |
| <i>B.</i> . . . . .       | buffet.                            | <i>Mar.</i> . . . . .  | Mardi.                 |
| <i>cent.</i> . . . . .    | centime.                           | <i>Mer.</i> . . . . .  | Mercredi.              |
| <i>ch.</i> . . . . .      | chaque.                            | <i>mil.</i> . . . . .  | millimètre.            |
| <i>cl.</i> . . . . .      | classe.                            | <i>min.</i> . . . . .  | minutes.               |
| <i>corres.</i> . . . . .  | corres-<br>pondance                | <i>mt.</i> . . . . .   | mont.                  |
| <i>déj.</i> . . . . .     | déjeuner.                          | <i>N.</i> . . . . .    | Nord.                  |
| <i>Dim.</i> . . . . .     | Dimanche.                          | <i>O.</i> . . . . .    | Onest.                 |
| <i>dr.</i> . . . . .      | droite.                            | <i>p.</i> . . . . .    | page.                  |
| <i>E.</i> . . . . .       | Est.                               | <i>s.</i> . . . . .    | siècle.                |
| <i>env.</i> . . . . .     | environ.                           | <i>S.</i> . . . . .    | Sud.                   |
| <i>fr.</i> . . . . .      | franc.                             | <i>Sam.</i> . . . . .  | Samedi.                |
| <i>g.</i> . . . . .       | gauche.                            | <i>S.</i> . . . . .    | Saint.                 |
| <i>h.</i> . . . . .       | heure.                             | <i>ser.</i> . . . . .  | service.               |
| <i>h. m.</i> . . . . .    | heures mi-<br>nutes.               | <i>St.</i> . . . . .   | Station.               |
| <i>hab.</i> . . . . .     | habitants.                         | <i>V.</i> . . . . .    | Ville.                 |
| <i>J.-S.-C.</i> . . . . . | Jonction-<br>Salonique<br>Cons/ple | <i>v.</i> . . . . .    | voir.                  |
| <i>Jeu.</i> . . . . .     | Jeu.                               | <i>Ven.</i> . . . . .  | Vendredi.              |
| <i>Kil.</i> . . . . .     | kilomètre.                         | <i>W.-L.</i> . . . . . | Wagon-Lit.             |
| <i>L.</i> . . . . .       | Ligne.                             | <i>W.-R.</i> . . . . . | Wagon-Res-<br>taurant. |
| <i>larg.</i> . . . . .    | largeur.                           | <i>†</i> . . . . .     | décédé.                |

**RÈGLES COMMUNES**

la sixième. Les élèves ont à peine appris les déclinaison régulières, qu'ils vont refaire méthodiquement pour cette ce qu'ils ont fait, sans réflexion, pour leur langue maternelle

Le Professeur est au tableau; en montrant sa tête, il l et écrit lentement, et tous ses élèves écrivent aussi :

|                  |  |
|------------------|--|
| <b>Tête,</b>     | <i>caput, pitis, (n.)</i> d'au capital, capitale, capitaux, capitaine, chap                          |
| <b>Cerveau,</b>  | <i>cerebrum, i, (n.)</i> d'où cérébral (fièvre cérébrale), etc.                                      |
| <b>Œil,</b>      | <i>oculus, i, (m.)</i> d'où oculiste, oculaire.  |
| <b>Front,</b>    | <i>frons, tis, (f.)</i> d'où frontal, fronton, frontispice, fronteau, etc.                           |
| <b>Nez,</b>      | <i>nasus, i, (m.)</i>    nasal, nasalité, naseau, nasillard, nasiller, nasille<br>nasarde, nasarder. |
| <b>Langue,</b>   | <i>lingua, e, (f.)</i>    lingual, linguiste, linguistique.  |
| <b>Dent,</b>     | <i>dens, tis, (m.)</i>    dentiste, dentition, dentier, dentaire, dental, de<br>telle, etc.          |
| <b>Barbe,</b>    | <i>barba, e, (f.)</i>    barbier, barbu, barbiche, barbifier, barbet, etc.                           |
| <b>Col,</b>      | <i>collum, i, (n.)</i>    collier, collet, collette, colleter.                                       |
| <b>Epaule,</b>   | <i>humerus, i, (m.)</i>    humerus, huméral.<br><i>scapula, arum, (f. pl.)</i>    scapulaire.        |
| <b>Côté,</b>     | <i>latus, teris, (n.)</i>    latéral, latéralement.  |
| <b>Poitrine,</b> | <i>pectus, oris, (n.)</i>    pectoral.   |
| <b>Estomac,</b>  | <i>stomachus, i, (m.)</i>    stomachique, stomacal.  |
| <b>Cœur,</b>     | <i>cor, cordis, (n.)</i>    cordial, cordialité, cordialement.                                       |
| <b>Sang,</b>     | <i>sanguis, inis, (m.)</i>    sanguin, sanguinaire, sangsue; saignant, sa                            |
| <b>Chair,</b>    | <i>caro, carnis (f.)</i>    carnavaul, carnassier, carnivore, carnation, etc                         |
| <b>Main,</b>     | <i>manus, us, (f.)</i>    manuserit, manoeuvre, manoeuvrer, manuel, ma<br>nutenlition.               |
| <b>Pied,</b>     | <i>pes, pedis, (m.)</i>    pédale.   |

N. B. On peut partager cette liste en deux leçons de six ou s chacune, et omettre les dérivés trop savants.

<http://gallica.bnf.fr/ark:/12148/bpt6k55584488/f19.image>

Glossaire, liste d'abréviations

<http://gallica.bnf.fr/ark:/12148/bpt6k55125514/f40.image>

Liste de définitions

| 1 <sup>er</sup> SEMESTRE (suite)                                    | 2 <sup>e</sup> SEMESTRE (suite)  |
|---|--|
| 2 <sup>e</sup> Explication et récitation d'auteurs.                 | 2 <sup>e</sup> Explication et récitation d'auteurs.                                      |
| 3 <sup>e</sup> Exercices oraux sur les textes expliqués.            | 3 <sup>e</sup> Exercices de conversation sur les textes expliqués.                       |
| 4 <sup>e</sup> Thèmes (Les reprendre de mémoire).                   | 4 <sup>e</sup> Thèmes d'imitation et d'application des règles.                           |
| 5 <sup>e</sup> <i>Grammaire</i> . — Révision du cours de troisième. | 5 <sup>e</sup> <i>Grammaire</i> . — Etude des verbes composés.                           |
| Syntaxe. — Les prépositions et les conjonctions.                    | Influence des préfixes et des particules sur la conjugaison et sur l'acception du verbe. |
| — Verbes irréguliers.   |  |
| — Récapitulation.   |  |

## CLASSE DE RHÉTORIQUE

| 1 <sup>er</sup> SEMESTRE  | 2 <sup>e</sup> SEMESTRE                                       |
|---|---|
| 1 <sup>er</sup> Lexicologie et exercices de conversation sur les mots appris.                   | 1 <sup>er</sup> Idiotismes et proverbes.                      |
| 2 <sup>e</sup> Explication et récitation d'auteurs.   | 2 <sup>e</sup> Formation et dérivation des mots.              |
| 3 <sup>e</sup> Lecture courante de morceaux faciles.  | 3 <sup>e</sup> Prosodie.                                      |
| 4 <sup>e</sup> Exercices de conversation sur les textes lus et expliqués.                       | 4 <sup>e</sup> Thèmes écrits et oraux.                        |
| 5 <sup>e</sup> Thèmes d'application des règles.   | 5 <sup>e</sup> Notions d'histoire littéraire sur les auteurs. |
| 6 <sup>e</sup> <i>Grammaire</i> . — Révision, en insistant sur les remarques et les exceptions. |   |

Comme le nouveau programme du baccalauréat parle d'un thème fait sans dictionnaire, de l'explication d'un texte et d'un entretien, il importe d'habituer de bonne heure les enfants à parler la langue vivante qu'ils ont choisie. Le maître devra donc, dès le commencement, non pas enseigner en anglais ou en allemand, mais dire quelques phrases très simples, qu'il traduira en français, si c'est nécessaire ; et il fera répéter. Tout le monde convient que pour parler une langue, il faut l'entendre parler, vivre en quelque sorte dans son milieu et s'exercer soi-même. Si les élèves sont ainsi forcés à reproduire, dès la quatrième, quelques phrases, ils s'habitueront peu à peu ; et, en troisième, le professeur pourra déjà se donner plus librement carrière. Dans ces phrases, le maître fera entrer surtout les mots déjà vus. Les enfants auront moins de peine à le comprendre et à répéter. Il sera utile et peut-être nécessaire

<http://gallica.bnf.fr/ark:/12148/bpt6k55125514/f80.image>

Texte composé sur deux colonnes



### ATTENTION

LES TABLES DES MATIÈRES ET INDEX, QUI ONT SOUVENT LA FORME D'UN TABLEAU ET QUI SONT TRANSCRIT DANS CERTAINS MARCHES NE SONT PAS CONCERNÉS PAR CETTE RÈGLE ; CES CONTENUS DOIVENT ÊTRE TRAITÉS SELON LES CONSIGNES DU « REFERENTIEL TABLES ». UN FICHER ALTO OCR BRUT EST CEPENDANT FOURNI POUR CES PAGES DANS LES DOCUMENTS NUMÉRIQUES.

- tous les autres types de tableaux (mise en page complexe, tableaux de texte et nombres, petit corps de police, etc.).

Ces contenus seront traités en **OCR brut**. Ils seront balisés :

**TAGREFS="TAG\_table TAG\_raw"**

#### 5.6.4 Encadrés

Il s'agit d'un encadré au sein d'une colonne (c'est-à-dire qui n'est pas un élément central.)

En règle générale, il n'y a pas un balisage particulier de l'encadré.

**+** EXEMPLE

UN ISLAM CRISPÉ

vement et exclusivement les préceptes supposés de l'islam primitif, en tant que système fermé, absolu et parfait, confondant société civile, société religieuse et société politique<sup>6</sup>. Pour le chef chiite, la décadence de la nation musulmane et du peuple iranien n'a d'autre raison que l'introduction de pouvoirs séculiers dans la société islamique, c'est-à-dire la laïcisation et l'ouverture de la société aux valeurs occidentales. Mais plus grave que l'idéologie religieuse de Khomeini lui-même est le soutien qu'elle a reçu des intellectuels en Occident tout comme en Orient, de même que la manipulation des médias à son profit, au détriment des autres composantes du soulèvement iranien. Exotisme raciste en Occident, aliénation culturelle et fascination du pouvoir en Orient auront à nouveau permis de « confisquer » à tout un peuple son soulèvement courageux contre la dictature et une fausse modernisation autoritaire.

**L'ÉGLISE CHÏTE ET LE POUVOIR \***

**M**énacé par la volonté totalitaire du régime, la classe religieuse se perçoit tenue en état de siège par les autorités laïques, qui elles-mêmes d'origine monarchique, libérale ou même islamique. Le chef de la fraction la plus puissante de l'Église dans la fin de l'expérience mousaddeghienne est, de ce point de vue, remarquable. Ce qu'elle exprime bien, Mousaddegh n'a pu empêcher de lui faire le rétablissement de l'absolutisme monarchique, au-delà des effets réels à sa politique à l'égard des grandes firmes privées, de la préservation des droits de la femme, à sa permission à l'égard du Toudéh qui l'empêche à la fois par l'usage d'espionnage et par son rôle politique afin d'appuyer aux vues des Américains car il est le seul rempart au communisme, aux valeurs qu'il représente avec l'Église nationale française de « l'équilibre nul ». C'est à lui-même, son désir de maintenir les religieux à l'écart du pouvoir. Pour lui qui, cependant, affirmait son attachement inébranlable aux principes d'égalité et de liberté, comme pour Mousaddegh, même devenu fermement politique, qui l'inspirait, le rôle historique de l'Église chiite est la critique du pouvoir - elle ne peut être à la fois le pouvoir et sa critique. L'appartenance de pouvoir de membres de clergé était par vous inspiré par l'idée qu'un gouvernement objet d'un large consensus et qui se place pas aux mains des religieux constituait un danger, une telle situation fait émettre des soupçons habituels entre l'Église et les classes populaires, la lutte d'un gouvernement impopulaire était dans l'ordre des choses non la faiblesse d'un gouvernement démocratique.

Paul Vialle

Paul Vialle, « Transmutation des aspects sociaux et révolution en Iran », in *Projetes révolutionnaires, juillet-septembre 1976*, pp. 31-35.

6. Cf. la traduction en arabe des conférences (francaises) Nadj en Irak par Hiram Khomani - *Al-Jihadoune Al-Islamiyye*, Dar el-Talib, Beyrouth, 1979. On les trouve avec l'introduction préface de la maison d'édition libanaise, républicaine raïssa d'Israël de Gharbi et qui présente la pensée de Khomeini comme une pensée religieuse traditionnelle et monarchique, et ce fait suscite de leur un rôle révolutionnaire dans le monde arabe.

34



5.6.5 Notes de bas de page

Les notes sont segmentées dans un bloc différent de celle du corps du texte. Chaque note doit être segmentée dans un bloc distinct (voir aussi section 5.5.4).

**+** EXEMPLE

Alors s'avança au-devant du chevalier sir Hugh le Héron, baron de Twisell et de Ford, gouverneur de Norham, qui le conduisit à la place d'honneur, au dais de l'estrade.

Le repas fut excellent et joyeux; et, pendant ce banquet, un ménestrel grossier du nord chanta sur la harpe le récit d'une sanglante inimitié; il dit comment — les farouches Thirwalls, tous les Riddleys, le robuste Willimondswick, Dick de Hardriding, Hughie de Hawdon, et Will o' the Wall, fondirent sur sir Albany Featherstonbough, et l'égorgeurent à Deadinan's Shaw (1).

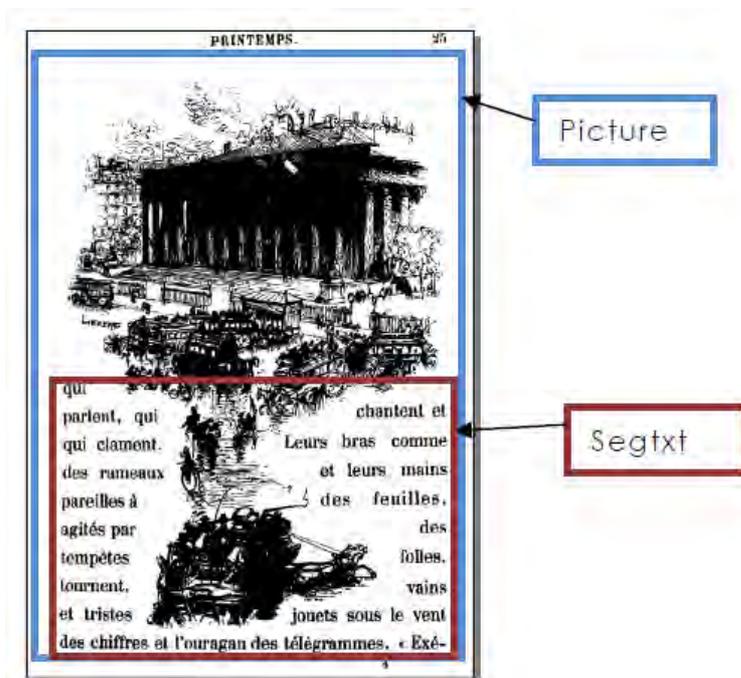
Marmion eut peine à écouter jusqu'au bout ce chant barbare; mais reconnaissant de la peine du ménestrel,

(1) Citation d'une vieille ballade chantée par les ménestrels, et qui est très-populaire en Écosse. — Éd.



### 5.6.6 Illustrations avec habillage de texte traversant

Il s'agit des cas où l'ordre de la lecture « traverse » ligne par ligne l'image, comme dans l'exemple ci-dessous.



Ce type de mise en page sera décrit avec un composedBlock comprenant des blocs textes pour le texte et un bloc Illustration autour de l'image, ainsi la lecture ligne à ligne de l'original est respectée (voir section 5.9.2).

### 5.6.7 Publicités et catalogues d'éditeur

Il s'agit de zones à caractère publicitaire qui présentent des éléments typographiques particuliers tels que polices curvilignes, logos, encadrés, etc. Elles seront traitées en OCR brut.



Ces pages sont de type « A » dans le manifeste numérique refNum.



EXEMPLE

AUX  
**PHARES**  
DE LA  
**BASTILLE**

|  |    |
|--|----|
| LE SUCCÈS superbe habille-<br>ment complet en<br>drap nouveauté.....                 | 26 |
| MARIAGE Habille-<br>ment complet, redingote, pan-<br>talon, gilet, le tout pour..... | 32 |
| L'INUSABLE magnifique<br>pardessus pour<br>hommes, drap nouveauté riche...           | 18 |
| 1 <sup>re</sup> COMMUNION Habille-<br>ment complet tout en drap noir fin.....        | 14 |

**SOLIDITÉ. ÉLÉGANCE. BON MARCHÉ**  
Ne se trouvent qu'aux  
**PHARES DE LA BASTILLE**  
5 et 7, place de la Bastille  
PARIS  
Envoi franco du MAGNIFIQUE CATALOGUE  
ILLUSTRÉ à toute personne qui en fait la  
demande.

J. H. ED. HEITZ, IMPRIMEUR-ÉDITEUR, STRASBOURG.  
Maison fondée en 1802 à Strasbourg. Dirigée depuis 1881 par les Heitz.

**BIBLIOTHECA ROMANICA.**  
DIRECTION | ED. SCHNEEGANS, Strasbourg  
PAUL HEITZ, Strasbourg

LISTE DES COLLABORATEURS:

|                             |                                |
|-----------------------------|--------------------------------|
| C. Appert, Bruxelles        | M. Leprieux, Berlin            |
| C. Baillet, Vienne          | F. Lottz, Düsseldorf           |
| F. Beck, Bamberg            | K. Michaux de Vescomtal, Paris |
| Aug. Becker, Leipzig        | P. Nati, Göttinge              |
| A. Closser, Chartres        | F. Neri, Strasbourg            |
| J. Clouet, Paris            | J. J. Orlitz, Vienne           |
| S. Dübendorf, Turin         | C. Ozando, Rome                |
| F. Düding, Bielefeld        | R. Palmarsch, Dрезно           |
| J. Fribolshagen, Strasbourg | A. Paris, Metz                 |
| Tu. Genard, Strasbourg      | P. Savi-Lopez, Naples          |
| G. Gogh, Vapiana            | R. Schenckhaus, Aachen         |
| H. Golub, Strasbourg        | Ed. Schneegans, Strasbourg     |
| T. G. Gruber, Strasbourg    | E. Suardi, Göttinge            |
| H. Haury, Paris             | L. Serrano, Catania            |
| E. Heiler, Heidelberg       | G. Tassin, Speyer              |
| Hopff, Strasbourg           | C. Thib, Strasbourg            |
| J. F. Halle, Berlin         | H. Yagow, Lyon                 |
| L. Jordan, Munich           | B. Weiss, Halle                |
| F. Kubler, Strasbourg       | W. Warbach, Vienne             |

Paris depuis 1881. — Le prix de chaque ouvrage est de 1 fr. 50.  
Chaque volume peut être fourni soit en belle toile et sur dent.  
Le prix de la reliure toile de R. 2 fr. 50.

**Bibliothèque française**

|  |  |
|--|--|
| Bilan. Eugénie Grandet. — Introduction par H. Giffet. 8082.  | Quatre. Mauvais de Journal, Lottz, Strasbourg. 182 10.   |
| Le Cabot des Antilles. — Int. par H. Giffet. 8083.   | La Héroïque Caracore. — Int. par F. Ed. Schneegans. 182 107.   |
| Boutemarchand. Le Barbier de Séville. — Int. par G. Gruber. 25 24.   | Lamarine Méditerranée. — Int. par F. Ed. Schneegans. 75 77.  |
| Bourgeois de Saint-Pierre, Paul et Virginie. — Int. par A. Paris. 117 118.   | La pulque Bourgeois. Poème satirique de Jan 1816. — Int. par M. Leprieux. 255.   |
| Bouffon. Art poétique. — Int. par E. Hopff. 94.  | Mars. Chronique. Passier. Wagons avec obstacles. — Int. par Th. Genard. 252 54.  |
| Le Lutin. — Int. par E. Hopff. 101.  | Molière. Le Misanthrope. — Int. par G. Gruber. 1.  |
| Châteaux populaires des XV <sup>e</sup> et XVII <sup>e</sup> siècles avec leurs mœurs. — Int. par Th. Genard. 180 182. | Les Femmes savantes. — Int. par G. Gruber. 2.  |
| Chateaubriand. Atala. — Int. par F. Ed. Schneegans. 84 85.   | L'Avant. — Int. par C. Thib. 46.   |
| — Noë. — Int. par F. Ed. Schneegans. 181.  | Terrils. — Int. par G. Gruber. 119.  |
| Contes de peuples politiques et satiriques du temps de la Fronde. — Int. par M. Leprieux. 257 258.                     | L'Étude des Français. — La critique de l'école des Français. — L'impression de Versailles. — Remerciement au roi. — Int. par F. Ed. Schneegans. 252 257. |
| Cornouille. Le Cid. — Int. par G. Gruber. 3.   | Le Balade imaginaire. — Int. par F. Heiler. 228 228.   |
| Horace. — Int. par C. Thib. 28.  | Les Fables. — Int. par F. Ed. Schneegans. 281.   |
| Ulysse. — Int. par C. Thib. 50.  | Le Bourgeois gentilhomme. — Int. par C. Thib. 249 250.   |
| Le Menace. — Int. par C. Thib. 92.   | Monsieur de Pourceaugnac. Int. par F. Ed. Schneegans. 250.   |
| Dissertation. Histoire de la méthode. — Int. par G. Gruber. 4.   | L'Amour malade. Int. par F. Ed. Schneegans. 250.   |
| Diderot. Le Paradoxe sur le Comédien. — Le Nivôse au Rhénan. — Int. par E. Lottz. 170 82.                              |  |

Si une zone de publicité comporte du texte et des éléments typographiques particuliers, l'ensemble de la zone est identifiée en publicité. Chaque bloc de texte contenu dans une publicité sera donc marqué comme un TextBlock avec une étiquette LayoutTag LABEL="advertisement" (cf. section 5.10).

Si plusieurs publicités sont adjacentes, elles ne sont pas groupées au sein d'un même bloc mais identifiées séparément.



NOTES

Si un bloc de publicité est composé sous forme d'un tableau, un double étiquetage LayoutTag "advertisement" et LayoutTag "table" sera utilisé.

Si une zone de publicité comporte du texte et des illustrations, les illustrations seront segmentées en bloc Illustration et ces blocs seront également étiquetés LayoutTag LABEL="advertisement".

**+** EXEMPLE



Si une zone de publicité comporte du texte et des illustrations mêlées, la totalité de la zone sera segmentée en bloc Illustration.

**+** EXEMPLE



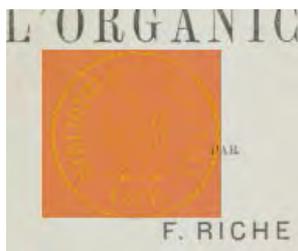
5.6.8 Texte sous tampon

Un bloc de texte placé sous un tampon sera marqué comme un TextBlock avec une étiquette LayoutTag LABEL="textStamped" (cf. section 5.10).

Le bloc de texte sera groupé avec le bloc tampon à l'aide d'un bloc composé.

Le texte contenu sera traité en OCR brut.

**+** **EXEMPLE**



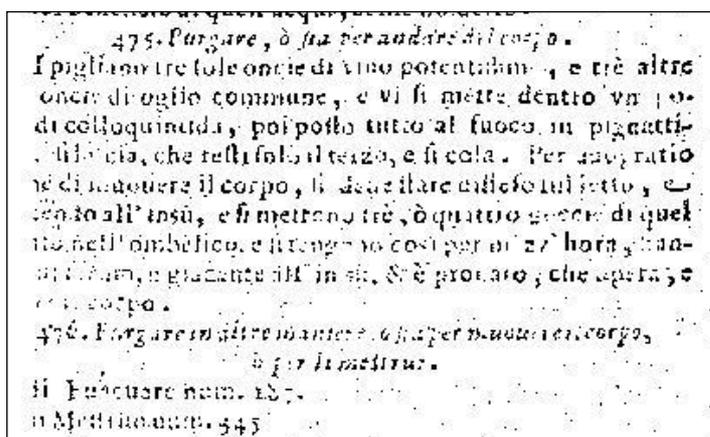
*Ici, le bloc "PAR"*

---

### 5.6.9 Texte illisible

Un bloc de texte illisible du fait de la dégradation physique du document, de problème de courbure ou de tout autre phénomène affectant la netteté de l'image, sera marqué en tant que TextBlock avec une étiquette LayoutTag LABEL = "illisible" (cf. section 5.10).

**+** **EXEMPLE**



---

La section 6.2 donne une définition précise de la notion de « texte illisible ».

**!** **ATTENTION**

LES ELEMENTS CONTENUS DANS UN BLOC ILLISIBLE (LIGNES ET MOTS) NE SERONT PAS EUX-MEMES ETIQUETES EN « ILLISIBLE ».

### 5.6.10 Texte manqué

Un élément textuel non reconnu en tant que tel par l'OCR sera segmenté en bloc texte de contenu "[texte manquant]", et identifié par une étiquette LayoutTag ayant pour valeur "missing" (cf. section 5.4).



Bloc avec texte en rotation non reconnu par l'OCR

### 5.6.11 Texte en OCR brut

Un bloc de texte traité en OCR brut sera marqué en tant que TextBlock avec une étiquette LayoutTag LABEL="raw". Ce typage est à réaliser uniquement :

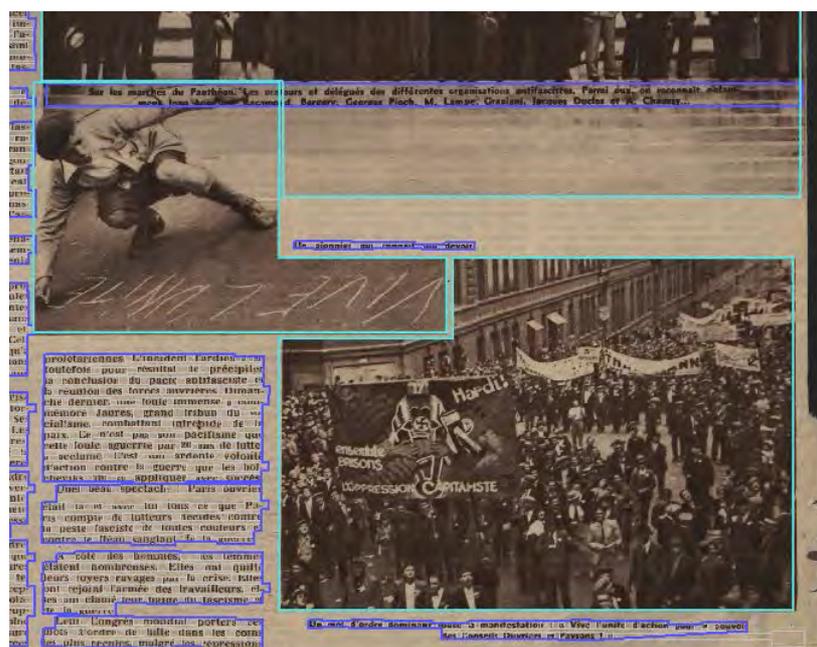
- pour certaines natures de contenus, en fonction des directives de la BnF,
- pour des prestations en qualité OCR garantie.

## 5.7 Illustrations

Tout élément de nature visuelle qui n'est à l'évidence pas une décoration, une lettrine ou un tampon sera traité sous la forme d'un bloc Illustration, sans chercher à statuer sur le lien éditorial qu'il a ou non avec le contenu textuel : cela couvre entre autres les illustrations, figures, diagrammes, photos, dessins, schémas, reproductions d'art, plans et cartes.

Par défaut, les illustrations seront segmentées sous forme de rectangle. Pour certains types documentaires, des formes géométriques autres que le rectangle (notamment le polygone) pourront être utilisées. A cet effet, on emploiera un élément ALTO Shape.

## EXEMPLE



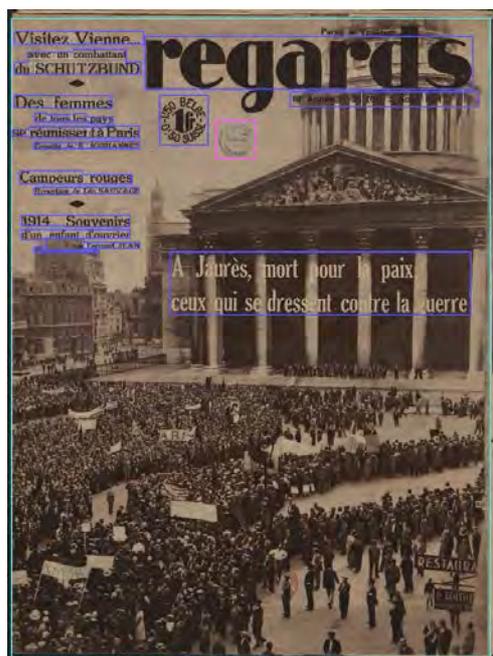
*Illustrations de presse polygonales*

Dans le cas d'utilisation de la seule forme rectangulaire, la couverture d'un bloc illustration ou élément graphique pourra être répartie en plusieurs segments rectangulaires, à condition que l'ensemble des blocs couvre la totalité de la zone illustrée.

## OBLIGATOIRE

On peut accepter que les illustrations combinées à d'autres éléments amènent à découper une zone « illustrée » en plusieurs blocs, à l'exception des partitions, formules et cartes qui doivent être découpées en un seul bloc.

En cas de mise en page complexe où textes et images sont mêlés, le bloc illustration aura l'emprise nécessaire et les blocs texte seront placés au-dessus.



### 5.7.1 Typage « illustration »

Le typage D (« dessin ») dans le fichier RefNum pourra être généré automatiquement à partir du typage ALTO : toute page dont le fichier ALTO contient un seul bloc Illustration (éventuellement accompagné d'un bloc de texte pour la légende de l'illustration) et qui n'a pas un autre typage prioritaire (P/E/T/I/R/L).



Pour générer un typage D, l'illustration doit occuper la majorité de la surface de la page.



Exemple de cas de typage D

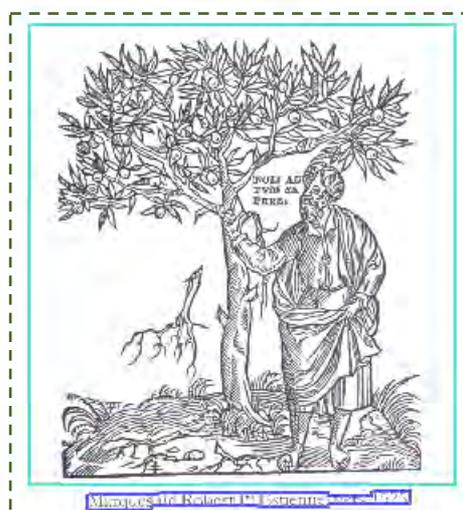
<http://gallica.bnf.fr/ark:/12148/bpt6k65164532/f19.image>

---

## 5.7.2 Imbrication de blocs Illustration et d'autres blocs, notamment des TextBlock

Quand une illustration s'accompagne d'une légende, il convient d'utiliser un ComposedBlock qui inclut l'image et la zone de texte.

La légende peut être placée sur l'illustration (cf. exemple du bas ci-dessous).





Bloc composé englobant illustration et légende



## NOTES

Selon les typologies de mise en page, cette règle ne sera pas toujours applicable et elle n'est pas soumise au taux qualité garantie.



## EXEMPLE



Maurice Ravel. *Finale de la Sonate pour piano.*  
Père Prellieur. *The modern music-master*, 1731.  
Gravure de J. Smith.

De la collection G. Thibault-de Cham-Dure :

- **Maurizio Cazzati**. *Riposta alle opposizioni fatte dal signor Giulio Cesare Arresi nella lettera al lettore posta nell'opera sua musicale*. Bologna, eredi del Dozza, 1663.



Cette plaquette, de soixante-douze pages, témoigne de la polémique entre deux compositeurs : Arresi, organisé à San Petronio de Bologne, et Cazzati (1616-1678), maître de chapelle de cette église, à propos de la technique utilisée dans une messe par ce dernier. Seul exemplaire connu en France. Ex-libris d'Henry Prunières.

- **Gottfried Keller**. *Rules or a complete method for attuning to play a thorough bass upon the harpsicord, organ or arclute... to which is added an exact scale for tuning the harpsicord or spinnet*. London, J. Walsh, s.d. (première moitié du XVIII<sup>e</sup> siècle).
- Six éditions de la méthode pour la basse continue de Keller, compositeur allemand du XVIII<sup>e</sup> siècle établi à Londres, paraurent dans cette ville de 1707 à 1728, outre celle-ci, encore inconnue et qui est la seule conservée en France.
- **Nicolaus Litanus**. *Rudiments musicae, in gratiam studiosae juvenutis diligenter comparata*. Augsburg, H. Steyner, 1536.
- Les cinquante-trois éditions de ce petit traité d'un maître de musique brandebourgeois publiées de 1533 à 1583 en Allemagne attestent sa grande popularité. La Bibliothèque nationale de France en possède cinq, de 1544 à 1557. Il fut d'abord préfacé par le réformateur Bugenhagen ; puis très augmenté à partir de 1557, son titre devint *Musica... nova regalis et exemplis adanata*. Ex-libris d'Alfred Corriot.
- **Père Prellieur**. *The modern music-master or the universal musician...* London, Printing Office, 1731.

De cet ouvrage de Prellieur, claveciniste

77

Mise en page de type presse, avec placement « souple » des légendes d'illustration

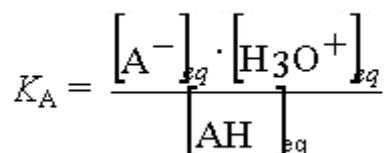
### 5.7.3 Formules chimiques, mathématiques

La structuration des formules est réalisée en mode image, avec un bloc Illustration doté d'une étiquette LayoutTag égale à "formula" (cf. section 5.10).

Il est permis de réunir plusieurs formules adjacentes (c'est-à-dire, quand il n'y a aucun autre bloc entre elles) en un seul bloc.

#### EXEMPLE

$$bO = AG - Am ; FO = FG - bm ;$$
$$bF = \sqrt{bO^2 + FO^2} ; \sin. bFO = \frac{bO}{bF}$$



#### NOTES

Quand des formules sont incluses dans des paragraphes de texte, elles ne sont pas traitées spécifiquement et la conversion OCR pourra donner des résultats illisibles.

### 5.7.4 Partitions

Les extraits de partition de musique sont encodés en bloc Illustration dotés d'un étiquette LayoutTag LABEL="musicScore" (cf. section 5.10).

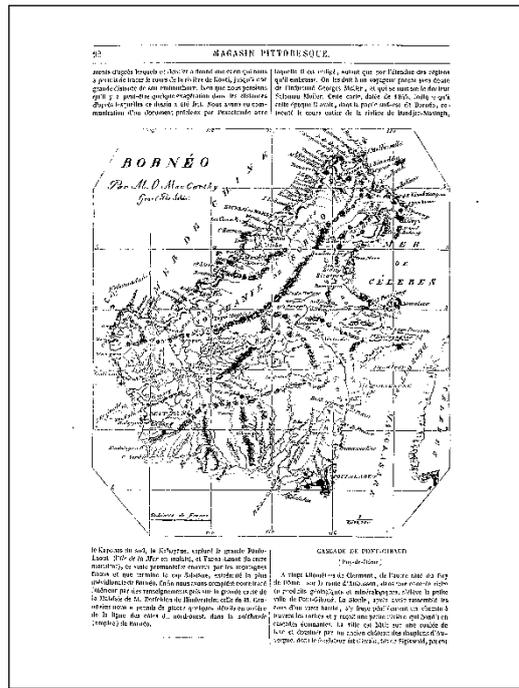
#### EXEMPLE



*Exemple de partition*

## 5.7.5 Cartes

Les illustrations reconnues comme étant des cartes géographiques seront typées spécifiquement en utilisant une étiquette LayoutTag égale à "map" (cf. section 5.10).



Exemple de carte

## 5.7.6 Alphabets non latins

Comme il a été dit section 4.1.2, les blocs de texte composés avec des alphabets non latins ou des systèmes d'idéogrammes seront décrits sous la forme de blocs image. Ces blocs seront des éléments Illustration dotés d'une étiquette LayoutTag LABEL="nonLatinScript" (cf. section 5.10).

## 5.7.7 Illustration en écriture manuscrite

Certaines illustrations sont la reproduction d'un texte manuscrit (extraits de manuscrit d'auteur, par exemple). Ces blocs seront décrits sous forme d'éléments Illustration (et non avec un élément graphique, cf. section 5.8.5).

l'humoriste appelait les deux sœurs « Une jolie paire de Cizos ».

Romieu remit une lettre pour son ex-collaborateur et ami Bayard, et toute la famille prit le chemin de Paris.

Engagée à l'essai au Gymnase, Rose débuta le 30 mars 1842 par le rôle d'Estelle, sous le nom de Marie C. Elle parut timide, plus correcte que gracieuse, mais disant bien. Au total, l'effet fut médiocre. M. Trubert, un marchand de rubans, qui dirigeait alors le Vaudeville, et Roqueplan, qui trônait aux Variétés, refusèrent de l'engager. Samson, à qui la jeune fille fut présentée, déclarait qu'il lui fallait dix-huit mois d'études. Les ressources de la famille s'épuisaient.

Chéri-Cizos fit une nouvelle tentative auprès de Monval, régisseur au Gymnase. On consentit à l'engager pour un an, à raison de 75 francs par mois, pour jouer les « en-cas ». Il a été vendu récemment (18 janvier 1902) une curieuse lettre de Monval à Rose Chéri. Cette lettre est du 5 juillet 1842. Le régisseur du Gymnase l'informe que M<sup>lle</sup> Natha-

rapide et brillante. Dès août 1842, Th. Gautier salue la nouvelle venue dans le *Premier chapitre* : « Enregistrons seulement les heureux débuts de M<sup>lle</sup> Rose Chéri, dont le nom charmant et le talent délicat ont favorablement disposé toute la critique. »

Il écrit encore le 13 septembre, à propos de *Céline* : « Cette jolie débutante (Rose Chéri) réussit beaucoup parce qu'elle est simplement une jeune fille toute naturelle, et n'a pas trop l'air d'une actrice; c'est le plus rare des talents. »

Le 16 décembre 1844, à propos de *Rebecca* : « La pièce est fort bien jouée par M<sup>lle</sup> Rose Chéri. »

Nous enregistrons : (les dates sont celles des feuilletons de Th. Gautier).

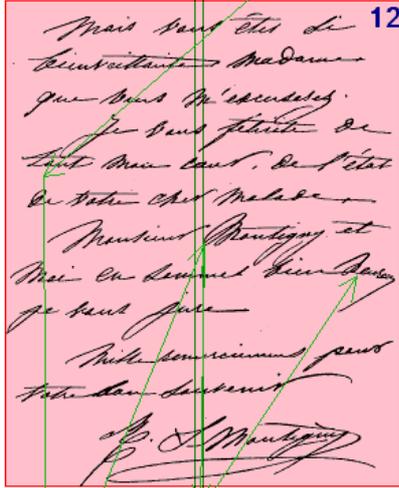
1845 janvier, M<sup>lle</sup> de Cérigny.

— juin, *Grande Dame et Grisette*.

(Assaut d'esprit entre M<sup>lle</sup> Rose Chéri et Désirée.)

1846 janvier, la *Loi Satique*.

« M<sup>lle</sup> Rose Chéri a joué son rôle avec ce naturel, cette intelligence et cette grâce qui lui assignent le premier rang parmi les actrices du vaudeville. »

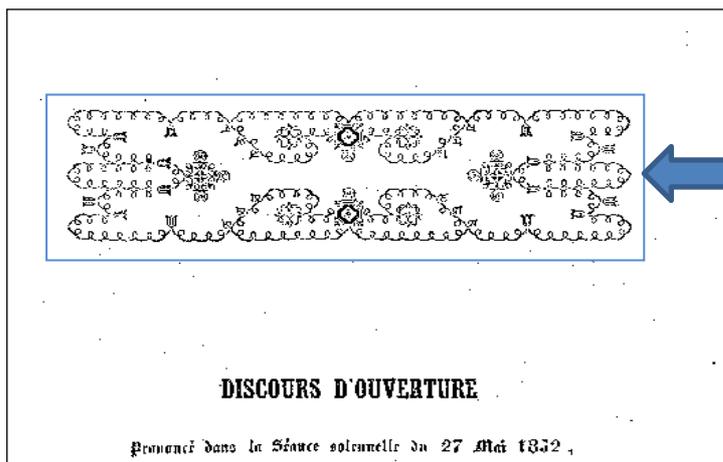


Exemple d'illustration sous forme d'écriture manuscrite

## 5.8 Éléments graphiques

### 5.8.1 Décorations et ornements

Les décorations et autres ornements sont à capturer en bloc GraphicalElement.



*Exemple de décoration*

---

### 5.8.2 Tampons

Les tampons sont à capturer en bloc GraphicalElement doté d'une étiquette LayoutTag LABEL="stamp" (cf. section 5.10).

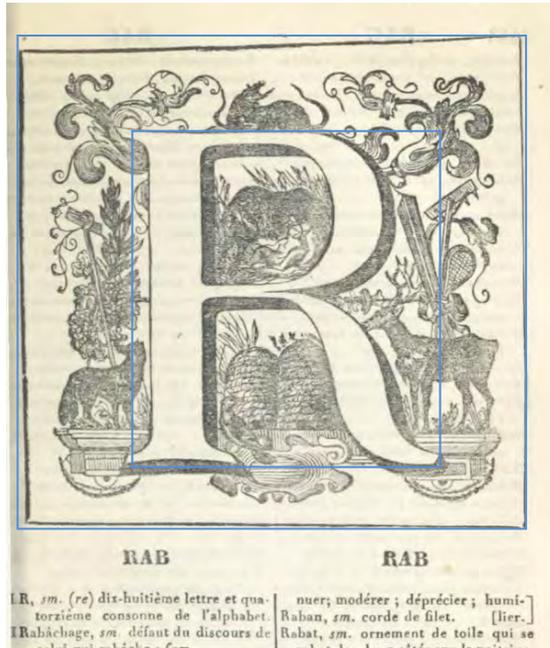


*Exemple de tampon*

---

### 5.8.3 Lettrines (lettres ornées)

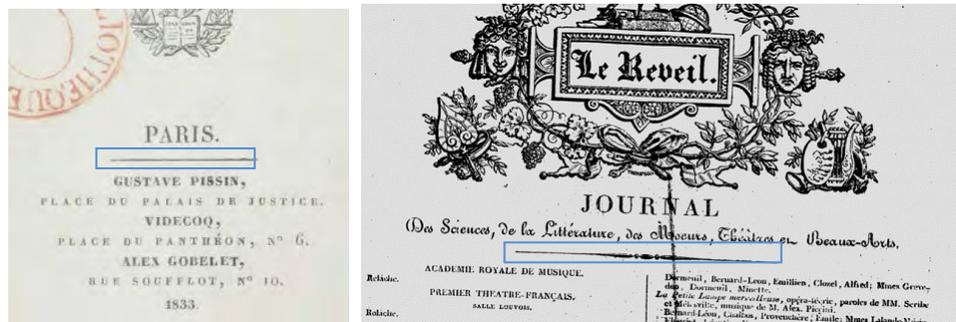
Les lettres ornées seront capturées en ComposedBlock avec un bloc texte pour le texte et un bloc GraphicalElement autour des ornements, doté d'une étiquette LayoutTag égale à "dropCap" (cf. section 5.10).

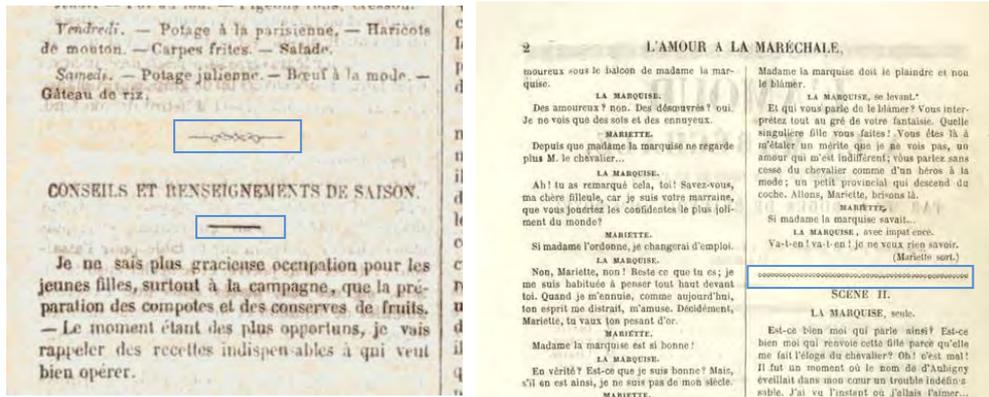


Exemple de lettrine (lettre R)

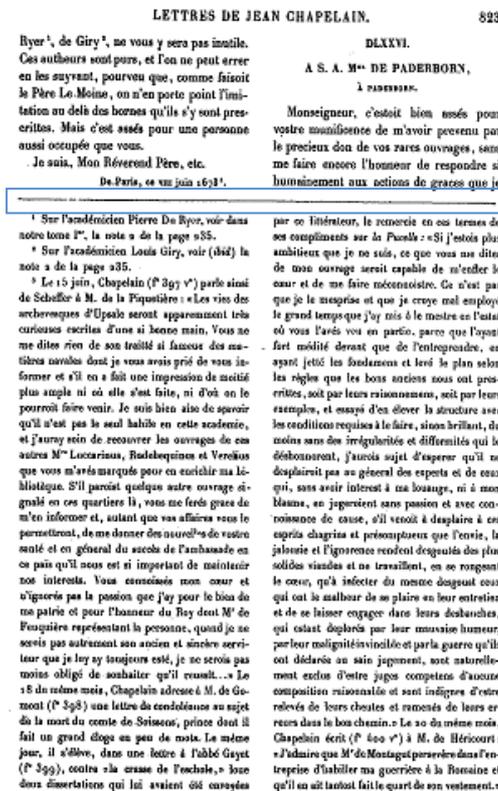
### 5.8.4 Traits de séparation

Les traits de séparation et les culs-de-lampe placés entre deux paragraphes doivent être décrits par un bloc GraphicalElement doté d'une étiquette LayoutTag LABEL="transition" (cf. section 5.10).



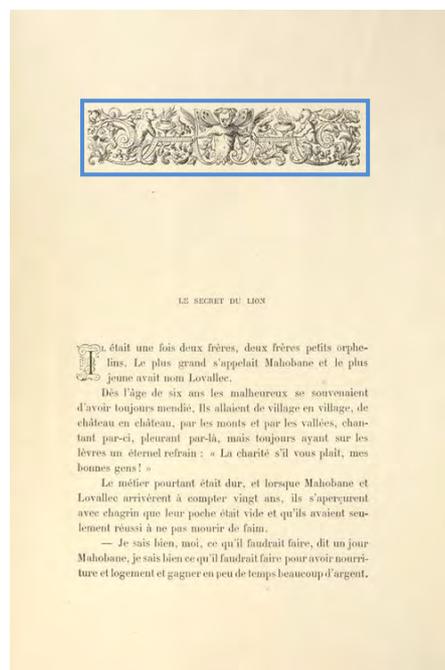
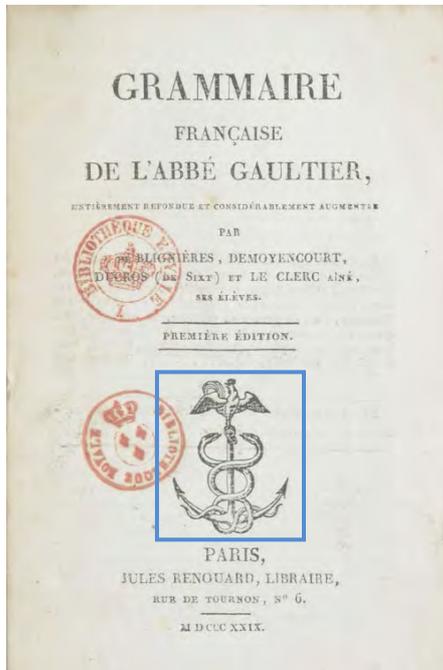
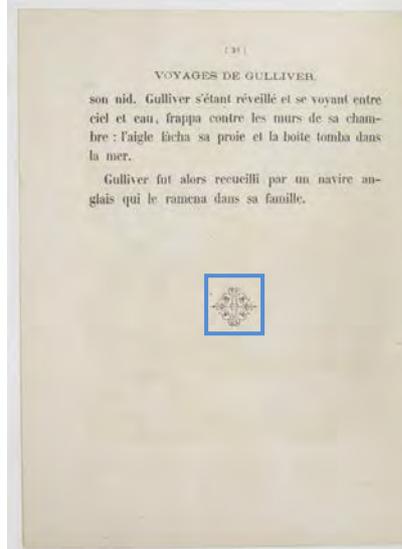


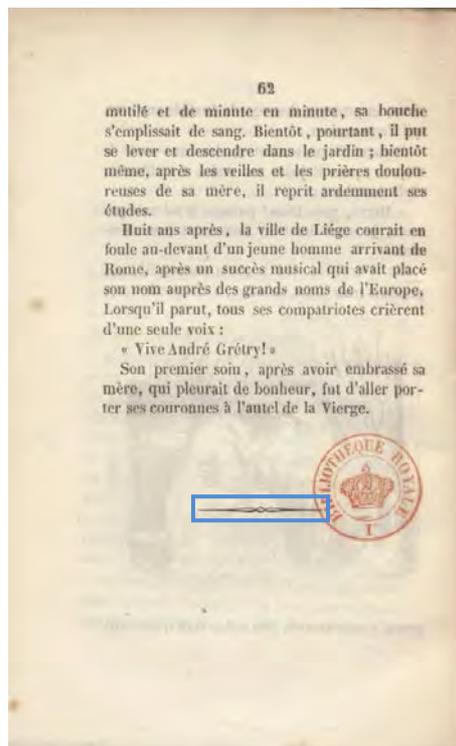
Les traits de séparation entre corps du texte et notes de bas de page sont à décrire en GraphicalElement avec une étiquette LayoutTag LABEL="footnoteSep".



Exemple de trait de séparation des notes de bas de page

Tous les autres types de culs-de-lampe (notamment ceux placés en fin de page ou en fin de chapitre) et de décorations sont à décrire en GraphicalElement (sans étiquette).





*Exemple de cul-de-lampe et décorations*

---

### *Contenu presse*

Les décorations et traits de séparation ornés seront segmentés.

Les traits de séparation non ornés ne seront segmentés que s'ils font toute la largeur de la page ou s'ils sont disposés à cheval sur au moins deux colonnes.

Les traits de séparation verticaux ne seront pas segmentés.



Exemples de traits de transition à segmenter

LES AUTRES TYPES DE SEPARATEURS DE TRANSITION POURRONT EVENTUELLEMENT ETRE SEGMENTES, MAIS EN AUCUN CAS OGERISES, NOTAMMENT CEUX POUVANT ETRE CONFONDUS AVEC DU TEXTE.

Le pas d'égalité sans en questions les priot dissipe dus. oc, a-t-il dit. te tous les me elle on- tous les de même du

2° Judas Furslemberg, né le 3 avril 1889, à Kiew (Russie), demeurant 42, rue d'Angoulême, dont la principale occupation et les ressources ordinaires étaient, d'après ses propres déclarations, de prendre des paris pour les champs de courses ;

(Voir la suite en 3<sup>e</sup> page)

**Le Ruban rouge**

**La promotion de l'instruction publique**

Demain paraîtra, à l'officiel, la promotion de l'instruction publique que nous avons annoncée depuis quelques jours et que nous donnons intégralement ci-dessous :

Né à Vienne, il est un-Mon- s'embrast pas ut beau- du haut sinciput is, tou- de. Seu- malheu- à il no-

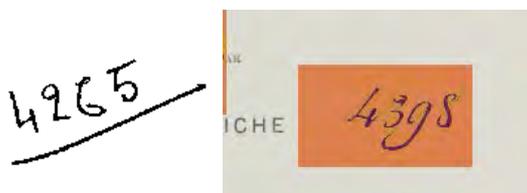
Peut-être un jour le D<sup>r</sup> Kornfeld nous apprendra-t-il que la rougeur des jeunes mariées, le soir de leurs noces, n'a pas d'autre cause. Avec les savants, il faut s'attendre à tout!

On rit volontiers de la maladie du sommeil lorsque l'on se trouve hors de ses atteintes. Quand on est près d'elle, on estime que c'est moins drôle. Les Belges ont en ce moment la

Plus joyeusement comme Les du Swift so glacé. Il qui, cha poussé v che vair picales. Au pe Wilde F ser tou se replie hérer le

### 5.8.5 Ecriture manuscrite

Les mentions manuscrites parfois présentes sur les pages des documents sont à capturer en GraphicalElement, avec une étiquette LayoutTag ayant pour valeur "manuscript" (cf. section 5.10).



*Exemples d'écriture manuscrite*

---



**CERTAINES PORTIONS D'ECRITURE MANUSCRITE SONT DES ILLUSTRATIONS. ELLES DOIVENT ETRE TRAITÉES SELON LES CONSIGNES ET EXEMPLE DE LA SECTION 5.7.**

### 5.8.6 Imbrication de blocs GraphicalElement et d'autres blocs, notamment des TextBlock

Les blocs GraphicalElement peuvent englober (via un ComposedBlock) d'autres blocs, notamment des TextBlock (voir section suivante).

## 5.9 Blocs composés

L'élément ComposedBlock est utilisé pour permettre le regroupement ou le recouvrement d'éléments de type bloc :

- regroupement logique (par exemple image et légende, cf. section 5.7.2),
- recouvrement physique : cf sections suivantes.

### 5.9.1 Texte au sein des illustrations ou des éléments graphiques

Ce cas concerne un texte pleinement intégré, c'est-à-dire que le texte fait partie intégrante de la zone graphique :

- de courts textes explicatifs au sein d'une illustration, notamment des légendes imbriquées dans le rectangle de l'illustration (cf. section 5.7.2).
- les textes au sein même d'un dessin, d'une œuvre d'art, d'un ornement, etc. (cf. section 5.8).

**+** EXEMPLE



Ce cas sera traité ainsi :

- le ou les TextBlock seront serrés autour du texte,
- l'image sera couverte par un bloc Illustration/GraphicalElement,
- les TextBlock et le bloc Illustration/GraphicalElement seront englobés par un ComposedBlock.

### 5.9.2 Imbrication d'illustrations ou d'éléments graphiques et de texte

Ce cas concerne du texte qui habille étroitement une illustration, en particulier, une illustration non rectangulaire.

**+** EXEMPLE



Ce cas sera traité ainsi :

- les TextBlock seront serrés autour du texte – jusqu'à un TextBlock distinct par ligne de texte,
- l'illustration sera couverte par un bloc Illustration, qui peut recouvrir entièrement ou partiellement les TextBlock ;

- l'ordre des blocs Illustration est libre, à condition que les TextBlock figurent dans l'ordre de la lecture. Les TextBlock peuvent être non contigus.

S'il y a recouvrement, les TextBlock et les GraphicalElement seront englobés par un ComposedBlock.

### 5.9.3 Ordre de lecture entre texte et illustrations ou éléments graphiques

Quand un ou plusieurs TextBlock partagent la largeur de la page ou de la colonne avec un ou plusieurs blocs Illustration/graphicalElement, seul compte l'ordre de lecture des TextBlock. L'ordre des blocs Illustration/graphicalElement n'est pas indiqué, ni l'emplacement de chacun d'entre eux avant, après ou au milieu des TextBlock.

## 5.10 Tableau récapitulatif de la structuration ALTO

Le tableau suivant détaille les différents cas d'usage des éléments structurants du format ALTO.

| TextBlock                                  |                                   |              |                             |
|--|-----------------------------------|--------------|-----------------------------|
| Cas  | Etiquette                         | Section      | Remarque                    |
| Paragraphe                                 |                                   | 5.6.1        |                             |
| Titre                                      | StructureTag : "titre1", "titre2" | 4.2.3, 5.6.2 |                             |
| Tableau                                    | LayoutTag : "table"               | 5.6.3        |                             |
| Encadré                                    |                                   | 5.6.4        |                             |
| Note (note de bas de page, note marginale) |                                   | 5.6.5        |                             |
| Texte habillé par une illustration         |                                   | 5.6.6        | utiliser un ComposedBlock   |
| Publicité                                  | LayoutTag : "advertisement"       | 5.6.7        |                             |
| Police "script"                            | LayoutTag : "scriptFonts"         | 4.2.3        |                             |
| Texte de l'ouvrage sous un tampon          | LayoutTag : "textStamped"         | 5.6.8        |                             |
| Texte illisible                            | LayoutTag : "illegible"           | 4.2.3, 5.6.9 |                             |
| Texte manquant                             | LayoutTag : "missing"             | 5.4          | texte non détecté par l'OCR |
| Texte en OCR brut                          | LayoutTag : "raw"                 | 5.6.10       | texte traité en OCR brut    |
| TextLine                                   |                                   |              |                             |
| Cas  | Etiquette                         | Section      | Remarque                    |
| Ligne illisible                            | LayoutTag : "illegible"           | 4.2.3, 5.6.9 |                             |

| String  |                              |              |                           |
|---|------------------------------|--------------|---------------------------|
| Cas   | Etiquette                    | Section      | Remarque                  |
| Mot illisible   | LayoutTag : "illegible"      | 6.3.2        |                           |
| Mot important   | OtherTag : "important"       | 6.2.1        |                           |
| Illustration  |                              |              |                           |
| Cas   | Etiquette                    | Section      | Remarque                  |
| Image, illustration, schéma, etc.                       |                              | 5.7          |                           |
| Formule mathématique ou chimique                        | LayoutTag : "formula"        | 5.7.1        |                           |
| Partition   | LayoutTag : "musicScore"     | 5.7.2        |                           |
| Carte, plan   | LayoutTag : "map"            | 5.7.3        |                           |
| Alphabet non latin                                      | LayoutTag : "nonLatinScript" | 4.1.2, 5.7.4 |                           |
| Illustration avec légende                               |                              | 5.7.6        | utiliser un ComposedBlock |
| GraphicalElement  |                              |              |                           |
| Cas   | Etiquette                    | Section      | Remarque                  |
| Décoration, cul de lampe                                |                              | 5.8.1        |                           |
| Tampon  | LayoutTag : "stamp"          | 5.8.2        |                           |
| Lettrine  | LayoutTag : "dropCap"        | 5.8.3        |                           |
| Trait de séparation entre paragraphes                   | LayoutTag : "transition"     | 5.8.4        |                           |
| Trait de séparation entre texte et notes de bas de page | LayoutTag : "footnoteSep"    | 5.8.4        |                           |
| Ecriture manuscrite                                     | LayoutTag : "manuscript"     | 5.8.5        |                           |

### 5.10.1 Etiquetage des éléments

On utilisera la syntaxe suivante :

```
<Tag ID="x" LABEL="étiquette" [DESCRIPTION="..."]/>
```

Une étiquette est associée à un élément à l'aide de l'attribut TAGREFS et d'un identifiant d'étiquette :

```
<Element ID="..." HPOS="..." VPOS="..." TAGREFS="id-x"/>
```

Une étiquette peut être *générique* (par exemple l'étiquette des « mots illisibles ») ou *instanciée*, c'est-à-dire associée à un élément particulier (par exemple l'étiquette de la « lettrine A du 2<sup>e</sup> bloc de la page 24 »).

Les identifiants des étiquettes génériques ne sont pas numérotés (la liste des étiquettes génériques est donnée à la fin de cette section) :

```
<LayoutTag ID="TAG_table" LABEL="table" />
```

Les identifiants des étiquettes instanciées sont numérotés (sur 3 digits, de 001 à 999) :

```
<LayoutTag ID="TAG_LT001" LABEL="table"  
DESCRIPTION="tableau de nombres" />
```



L'attribut `DESCRIPTION` est optionnel. Lorsqu'il est présent, il informe du contenu concerné par l'étiquette (soit en le décrivant, soit par son contenu textuel lui-même).

Une étiquette doit référencer le plus haut niveau disponible, par exemple un bloc et non tous les mots du bloc.

Un contenu à cheval sur deux pages peut être étiqueté par la même étiquette.

Un contenu peut être étiqueté par plusieurs étiquettes, en séparant les identifiants d'étiquette par un espace : `TAGREFS="id1 id2"`

### *Layout*

Les étiquettes `LayoutTag` servent à typer la nature des éléments de contenu (tableau, publicité, formule mathématique, carte, etc.).



Étiquettes génériques :

```
<LayoutTag ID="TAG_table" LABEL="table" />
```

```
<LayoutTag ID="TAG_textStamped" LABEL="textStamped" />
```

```
<LayoutTag ID="TAG_manuscript" LABEL="manuscript" />
```

```
<LayoutTag ID="TAG_scriptFonts" LABEL="scriptFonts" />
```

```
<LayoutTag ID="TAG_footnoteSep" LABEL="footnoteSep" />
```

Étiquettes instanciées :

```
<LayoutTag ID="TAG_LT001" LABEL="table" DESCRIPTION="tableau de nombres" />
```

```
<LayoutTag ID="TAG_LT008" LABEL="formula" DESCRIPTION="4 formules  
d'algèbre" />
```

```
<LayoutTag ID="TAG_LT009" LABEL="dropCap " DESCRIPTION="A" />
```

```
<LayoutTag ID="TAG_LT010" LABEL="nonLatinFont " DESCRIPTION="idéogrammes chinois"/>
```

### *Structure*

Les étiquettes StructureTag servent à décrire la structure logique des contenus. On utilisera la syntaxe suivante :

```
<StructureTag ID="x" TYPE="Structural" LABEL="étiquette" [DESCRIPTION="..."]/>
```



**Étiquettes génériques :**

```
<StructureTag ID="TAG_title1" TYPE="Structural" LABEL="title1"/>
```

```
<StructureTag ID="TAG_title2" TYPE="Structural" LABEL="title2"/>
```

**Étiquettes instanciées :**

```
<StructureTag ID="TAG_ST008" TYPE="Structural" LABEL="title1" DESCRIPTION="Chroniques littéraires"/>
```

### *Autres cas*

Les étiquettes OtherTag seront utilisées pour décrire les mots importants (cf. section 6.2.1) :

```
<OtherTag ID="x" [TYPE="..."] LABEL="étiquette" />
```

L'attribut TYPE pourra être utilisé pour préciser les types de mot important, par exemple les entités nommées (EN) ou les mots présents dans les titres.



**Étiquettes génériques :**

```
<OtherTag ID="TAG_important" LABEL="important"/>
```

```
<OtherTag ID="TAG_importantEN" TYPE="EN" LABEL="important"/>
```

```
<OtherTag ID="TAG_importantTL" TYPE="titre" LABEL="important"/>
```

**Étiquettes instanciées :**

```
<OtherTag ID="TAG_OT001" TYPE="EN" LABEL="important" DESCRIPTION="Charles Dickens"/>
```

## Liste des étiquettes génériques

| TextBlock                         |  |              |
|-----------------------------------|--|--------------|
| Cas                               | Etiquette  | Section      |
| Titre                             | <code>&lt;Structure Tag ID="TAG_title1" TYPE="Structural" LABEL="title1" /&gt;</code><br><code>&lt;Structure Tag ID="TAG_title2" TYPE="Structural" LABEL="title2" /&gt;</code><br><code>&lt;Structure Tag ID="TAG_title3" TYPE="Structural" LABEL="title3" /&gt;</code><br><code>&lt;Structure Tag ID="TAG_title4" TYPE="Structural" LABEL="title4" /&gt;</code> | 4.2.2, 5.6.2 |
| Tableau                           | <code>&lt;LayoutTag ID="TAG_table" LABEL="table" /&gt;</code>  | 5.6.3        |
| Publicité                         | <code>&lt;LayoutTag ID="TAG_advertisement" LABEL="advertisement" /&gt;</code>  | 5.6.7        |
| Police script                     | <code>&lt;LayoutTag ID="TAG_scriptFonts" LABEL="scriptFonts" /&gt;</code>  | 4.2.3        |
| Texte de l'ouvrage sous un tampon | <code>&lt;LayoutTag ID="TAG_textStamped" LABEL="textStamped" /&gt;</code>  | 5.6.8        |
| Texte illisible                   | <code>&lt;LayoutTag ID="TAG_illegible" LABEL="illegible" /&gt;</code>  | 4.2.3, 5.6.9 |
| Texte manqué                      | <code>&lt;LayoutTag ID="TAG_missing" LABEL="missing" /&gt;</code>  | 5.4          |
| Texte en OCR brut                 | <code>&lt;LayoutTag ID="TAG_raw" LABEL="raw" /&gt;</code>  | 5.6.10       |
| TextLine                          |  |              |
| Cas                               | Etiquette  | Section      |
| Ligne illisible                   | <code>&lt;LayoutTag ID="TAG_illegible" LABEL="illegible" /&gt;</code>  | 4.2.3, 5.6.9 |
| String                            |  |              |
| Cas                               | Etiquette  | Section      |
| Mot illisible                     | <code>&lt;LayoutTag ID="TAG_illegible" LABEL="illegible" /&gt;</code>  | 6.3.2        |
| Mot important                     | <code>&lt;OtherTag ID="TAG_important" LABEL="important" /&gt;</code>   | 6.2.1        |
| Illustration                      |  |              |
| Cas                               | Etiquette  | Section      |
| Formule mathématique ou chimique  | <code>&lt;LayoutTag ID="TAG_formula" LABEL="formula" /&gt;</code>  | 5.7.1        |
| Partition                         | <code>&lt;LayoutTag ID="TAG_musicScore" LABEL="musicScore" /&gt;</code>  | 5.7.2        |
| Carte, plan                       | <code>&lt;LayoutTag ID="TAG_map" LABEL="map" /&gt;</code>  | 5.7.3        |
| Alphabet non latin                | <code>&lt;LayoutTag ID="TAG_nonLatinScript" LABEL="nonLatinScript" /&gt;</code>  | 4.1.2, 5.7.4 |
| GraphicalElement                  |  |              |
| Cas                               | Etiquette  | Section      |

|   |  |       |
|---|--|-------|
| Tampon  | <LayoutTag ID="TAG_stamp" LABEL="stamp" />             | 5.8.2 |
| Lettrine  | <LayoutTag ID="TAG_dropCap" LABEL="dropCap" />         | 5.8.3 |
| Trait de séparation<br>entre paragraphes                      | <LayoutTag ID="TAG_transition" LABEL="transition" />   | 5.8.4 |
| Trait de séparation<br>entre texte et notes<br>de bas de page | <LayoutTag ID="TAG_footnoteSep" LABEL="footnoteSep" /> | 5.8.4 |
| Écriture manuscrite   | <LayoutTag ID="TAG_manuscript" LABEL="manuscript" />   | 5.8.5 |

## 6. QUALITE DE LA RECONNAISSANCE OCR

---

### 6.1 Qualité de la transcription du texte

La montée en qualité du texte, lorsqu'elle est demandée sur une prestation via un taux OCR qualité garantie ou haute qualité, a pour objectif d'atteindre un certain taux de reconnaissance des mots composant le document d'origine.



Ce taux est mesuré avec pour granularité le mot (élément ALTO <String>). Il est renseigné à l'échelle de chaque page ainsi qu'à l'échelle du document. C'est la valeur au document numérique qui fait foi sur les marchés.

Les mots reconnus par le moteur OCR sont marqués avec un critère de confiance calculé par le moteur OCR, selon le degré de reconnaissance du mot, en utilisant l'attribut wc de l'élément String.

Ce critère de confiance varie de 0 (mot non reconnu) à 0,99 (mot reconnu sans doute).



Chaque mot corrigé par un opérateur sera identifié dans le contenu ALTO (par exemple en positionnant à 1,0 le critère de confiance WC).

#### 6.1.1 Correction ciblée

Pour arriver au taux qualité attendu, le prestataire devra traiter prioritairement les mots dit « importants ». Ces mots importants se caractérisent par des caractéristiques typographiques ou intellectuelles, notamment :

- les mots avec une majuscule à l'initiale,
- les mots composés en majuscules,
- les éléments en gras ou en italique,
- les mots présents dans les titres,
- les entités nommées : noms propres, noms de lieu, noms de personne ou d'institution, etc.



Chaque mot identifié comme important sera repéré dans le contenu ALTO à l'aide d'une étiquette OtherTag de valeur "important" (cf. section 5.10).

## 6.2 Qualité de la segmentation

Le résultat de la segmentation doit permettre de faire correspondre le texte issu de l'OCR à l'image par transparence grâce au calcul des coordonnées de la position des éléments dans l'image.

Pour atteindre le taux qualité attendu, il est nécessaire de corriger (structuration, typage, ...) le résultat proposé par le moteur OCR.



La qualité de segmentation est mesurée à l'échelle du document numérique (ensemble des pages).

## 6.3 Déqualification de contenus dans un document

Il s'agit d'une opération d'identification des blocs de texte (TextBlock) ou des mots (String) difficilement océrisable, du fait de leurs caractéristiques propres. Ces zones de texte sont donc exclues de fait du périmètre du taux qualité garantie et elles sont traitées en OCR brut.

Cette identification se fait selon deux axes :

- par types de contenus,
- par lisibilité des contenus.

### 6.3.1 Déqualification par types de contenu

Cette déqualification s'appuie sur les résultats de la phase de segmentation/structuration, qui a conduit à typer certains contenus de manière objective, selon leur nature physique ou logique.

Ces critères discriminants quant aux types des contenus sont les suivants :

- textes dont la langue majoritaire n'est pas le français (cf. section 4.1.3),
- textes composés en polices manuscrites (cf. section 4.2.3),
- textes composés dans des polices non reconnues par le moteur OCR (cf. section 4.2.3),
- tableaux à traiter en OCR brut (cf. section 5.6.4),
- publicités (cf. section 5.6.8) : il s'agit de zones à caractère publicitaire qui présentent des éléments typographiques particuliers tels que polices curvilignes, logos, encadrés, etc.
- textes placés sous un tampon (cf. section 5.6.9).

Cette typologie peut faire l'objet d'un référentiel commun d'exemples type, qui sera alimenté d'un commun accord au fur et à mesure des cas spécifiques rencontrés.

Dans le cadre de l'amélioration continue des processus de production, cette opération pourra faire l'objet d'un processus plus automatisé. Dans ce cas, l'évolution du processus sera soumise à l'approbation de la BnF.

### 6.3.2 Déqualification des mots ou blocs illisibles

Les zones de texte illisibles (cf. section 5.6.10) sont définies non du fait de leur type logique mais à partir de leur nature physique :

- zones près de la reliure de l'ouvrage et présentant une forte courbure de l'image,
- zones affectées par une dégradation du support physique de l'œuvre, laquelle affecte la netteté, le contraste ou même l'intégrité des contenus
- zones affectées par un problème de transparence, de migration d'encre, etc.

#### *Mots*

Un mot est considéré comme illisible si un opérateur humain ne peut le déchiffrer à l'œil nu, ou s'il ne parvient à le faire qu'avec un fort degré de supposition et d'interprétation du contexte du texte.

La lisibilité est évaluée à partir de l'image binarisée ou de l'image source en cas de doute.

La déqualification opérée à l'échelle d'un mot consiste à typer le mot concerné (élément String) en lui affectant une étiquette "illegible" (cf. section 5.10).

#### *Lignes ou blocs*

L'illisibilité d'une ligne ou d'un bloc est définie par une valeur seuil de la proportion de mots illisibles dans la ligne ou le bloc. Ce seuil est défini à 50 %.

La déqualification opérée à l'échelle d'une ligne ou d'un bloc de texte consiste à typer l'élément concerné (TextLine ou TextBlock) en lui donnant une étiquette "illegible" (cf. section 5.10).



Les mots inclus dans un bloc ou une ligne illisible n'ont pas à être étiquetés "illegible".

### 6.3.3 Limites au principe de déqualification

Un document demandé par la BnF en taux qualité garantie et dont plus de 50 % du contenu textuel (relativement au nombre de mots) est déqualifié (selon les critères décrits aux sections 6.3.1 et 6.3.2) entre automatiquement dans le champ d'action du principe de déqualification du taux qualité, tel que décrit à la section 3.3.3. Il est alors livré par le prestataire en OCR brut.

## 6.4 Déqualification du taux qualité sur un document

Des demandes de déqualification en OCR brut peuvent être demandées (modalités à préciser pour chaque marché) a priori par le prestataire pour des documents dont les qualités physiques ne sont pas suffisantes pour atteindre la qualité garantie :

- les documents tâchés, bruités ou maculés (ex. des microfiches, microfilms) ;
- les documents à contraste insuffisant (fond foncé/tramé avec superposition de texte, ex : livre d'enfants) ;

- les documents avec forte transparence (ex : presse, dictionnaire, ouvrage à papier pelure) ;
- les documents dont la langue majoritaire n'est pas en taux garanti (latin, etc.).



Ces demandes de déqualification se font avant tout prétraitement OCR, sans utiliser le mécanisme de déqualification des contenus exposé à la section 6.3.

## 6.5 Déqualification ou refus de documents

Les documents connus pour être quasi intégralement inexploitable en OCR (documents manuscrits, composés en polices Fraktur ou gothique, avec un alphabet non latin, etc.) doivent être refusés ou déqualifiés en OCR brut par le prestataire pour la prestation OCR en fonction du type de marché.

Ces documents se reconnaissent à ce qu'ils se transcrivent sous la forme de blocs de texte vides ou très fortement bruités.

## 7. CONTROLE DE LA QUALITE

---

Le contrôle de la qualité est assuré par plusieurs moyens :

- des contrôles automatiques appliqués sur les contenus au format ALTO,
- un contrôle par échantillonnage visuel.

Au terme du contrôle, la BnF prononce le rejet ou l'acceptation des documents livrés par le prestataire.

Cette section décrit les critères de rejets ou d'acceptation

### 7.1 Contrôle automatique ALTO

Un contrôle automatique exhaustif de format est appliqué sur tous les fichiers ALTO avant le passage en contrôle par échantillonnage visuel. Ce contrôle émet des erreurs ainsi que des avertissements.

Une erreur entraîne le rejet du document.



SI CE CONTROLE EMET UNE ERREUR SUR UNE PAGE ALTO D'UN DOCUMENT, L'ENSEMBLE DES PAGES ALTO DU DOCUMENT SONT ECARTEES ET NE PASSENT PAS EN CONTROLE PAR ECHANTILLONNAGE VISUEL. LE DOCUMENT EST DONC REJETE DES CETTE ETAPE.

Ces contrôles automatiques sont de plusieurs natures :

- Nommage du fichier ALTO.
- Codage du fichier ALTO : UNICODE UTF-8 sans BOM (*Byte Order Mark*).
- Validation des fichiers ALTO relativement au schéma XML ALTO.
- Présence et format des identifiants et des différents attributs.
- Positionnement des blocs relativement aux dimensions de la page.
- Chevauchement entre blocs :
  - entre éléments de même niveau (ex: entre deux éléments String) ;
  - entre éléments de haut niveau (ex: entre un TopMargin et un PrintSpace) ;
  - entre éléments enfants et parents (ex: entre un TextLine et un TextBlock).

La tolérance est par défaut de 5%. Elle pourra être adaptée selon les prestations.

- Format des éléments de production : opérations, agents, résultats (schéma XML detailsOperation.xsd).



Les modalités de ce contrôle seront détaillées dans une charte de contrôle OCR/ALTO élaborée conjointement par la BnF et le prestataire.

## 7.2 Contrôle par échantillonnage visuel

Le contrôle par échantillonnage visuel opère sur des lots de documents constitués selon un plan d'échantillonnage adapté à chaque marché. Il vise à contrôler deux aspects du traitement OCR :

- la qualité de la reconnaissance du texte,
- la qualité de la segmentation/structuration.



LA DETECTION D'UNE NON-CONFORMITE SUR UN DOCUMENT ENTRAINE LE REJET DES DOCUMENTS COMPOSANT LE LOT DE CONTROLE.

LA DETECTION D'UNE NON-CONFORMITE STRUCTURELLE SUR TOUS LES DOCUMENTS COMPOSANT LE LOT DE CONTROLE ENTRAINE UN AUDIT DU PROCESSUS DE PRODUCTION ET EVENTUELLEMENT LE REJET DE TOUS LES DOCUMENTS DEJA PRODUITS.

### 7.2.1 Qualité de la reconnaissance du texte

La qualité de la reconnaissance du texte est mesurée au mot.

Une erreur de reconnaissance correspond à tout mot transcrit erronément par rapport à l'image d'origine. Un mot est considéré comme erroné quel que soit le nombre de signes erronés qu'il contient.

Pour un document, la qualité se calcule sur la population des mots présents dans toutes les pages du document, en tenant compte des règles définies pour le taux OCR considéré (brut, qualité garantie, qualité éditoriale, cf. section 4.1.3) :

- prise en compte des parties déqualifiées,
- jeu de caractères.



Un document peut avoir des parties dont le taux de reconnaissance est supérieur ou inférieur à la qualité attendue, mais l'ensemble du document doit valider la qualité attendue.

#### *Phase de test*

En numérisation de masse, la qualité de l'OCR est calculée à l'aide d'un taux de confiance estimé (en l'absence de vérité terrain). Ce taux (en %) est donnée par la formule :

$$\text{taux de confiance} = \frac{\sum (\text{wc des mots en qualité garantie})}{\text{cardinal des mots en qualité garantie}}$$

Le taux de qualité réel (spécifié pour chaque marché de numérisation comme qualité garantie à atteindre) est donnée par la formule :

$$\text{taux réel} = \frac{\text{cardinal des mots justes parmi les mots en qualité garantie}}{\text{cardinal des mots en qualité garantie}}$$

Durant la phase de test, on évaluera l'écart moyen entre la qualité estimée (taux de confiance) et la qualité réelle (taux réel). Cette évaluation sera réalisée sur des documents proches du taux qualité garantie du marché.

Cet écart sera utilisé en phase de production pour étalonner la valeur du taux de confiance (et donc l'effort de correction manuelle du texte nécessaire pour atteindre la qualité garantie) :

$$\text{taux de confiance} \times \text{coefficient d'étalonnage} \geq \text{taux garanti}$$



### OBLIGATOIRE

Les informations présentes dans les fichiers ALTO (cf. section 7.2.3) doivent permettre à la BnF de recalculer automatiquement les taux qualité fournis par le prestataire.

#### *Phase de production*

Le contrôle par échantillonnage visuel opère sur des lots de documents constitués selon un plan d'échantillonnage.

Le contrôle de la qualité de l'OCR porte sur l'écart entre la qualité attendue et la qualité constatée par la BnF.



### NOTES

Les modalités de ce contrôle seront détaillées dans une charte de contrôle OCR/ALTO élaborée conjointement par la BnF et le prestataire.

## 7.2.2 Qualité de la segmentation/structuration

Le contrôle de la segmentation/structuration s'applique à contrôler dans une page ALTO les éléments PrintSpace, xxxMargin, TextBlock, TextLine String, Illustration, GraphicalElement et ComposedBlock.

La granularité du contrôle est la page. Le niveau de qualité de segmentation/structuration acceptable est spécifique à chaque marché (cf. CCTP). Il indique donc le nombre maximum de pages non conformes dans un document. Pour les documents à faible pagination (presse), la mesure sera réalisée sur le nombre total de pages de l'échantillon de contrôle.



### NOTES

Ce contrôle ne s'applique pas sur les documents traités en OCR brut.

Pour une page à contrôler, le contrôle de la segmentation/structuration se décompose en plusieurs contrôles réalisés selon des modalités spécifiques à chaque contrôle :

| Contrôle                     | Modalité           | Résultat   |
|------------------------------|--------------------|--|
| Typage des pages             | contrôle unitaire  | correct/incorrect                                  |
| Identification du PrintSpace | contrôle unitaire  | correct/incorrect                                  |
| Identification des xxxMargin | contrôle unitaire  | correct/incorrect                                  |
| Ordre de lecture             | contrôle unitaire  | correct/incorrect                                  |
| Typage des blocs             | contrôle numérique | taux d'erreur $\leq$ (1 - taux qualité garantie) ? |
| Segmentation des blocs       | contrôle numérique | taux d'erreur $\leq$ (1 - taux qualité garantie) ? |



CHACUN DE CES CONTROLES DOIT ETRE SATISFAIT, CE QUI IMPLIQUE :

- QU'UN SEUL CONTROLE UNITAIRE EN ECHEC INDUIT LA NON-CONFORMITE DE LA PAGE
- QU'UN SEUL CONTROLE NUMERIQUE EN ECHEC INDUIT LA NON-CONFORMITE DE LA PAGE.

#### *Typage des pages*

Le contrôle vérifie le typage des pages et de leur orientation (cf. sections 5.1 et 5.2).

#### *Identification du PrintSpace et des xxxMargin*

Le contrôle vérifie l'identification du PrintSpace et des xxxMargin (cf. section 5.3).

#### *Ordre de lecture*

Le contrôle vérifie que l'ordre de lecture logique de la page est bien respecté (cf. section 5.5).

#### *Typage des blocs*

Le contrôle vérifie le typage des blocs identifiés dans la page (cf. sections 5.6 à 5.8) :

- nature des blocs segmentés (par ex. confusion blocs texte/bloc illustration),
- typage des blocs (par ex. confusion blocs de texte/blocs de publicité),
- blocs déqualifiés (par ex. bloc illisible abusif).

Le taux d'erreur est calculé ainsi :

$$\text{taux d'erreur typage} = \sum \text{erreurs typage} / \text{nombre de blocs}$$

Il doit être inférieur au taux qualité attendu pour la segmentation (cf. CCTP).

#### *Segmentation des blocs*

Le contrôle vérifie la bonne segmentation des blocs :

- taille et position des blocs relativement aux contenus (sauf pour les blocs déqualifiés en illisible),

- découpage cohérent des paragraphes du contenu textuel en blocs de texte (sauf pour les blocs déqualifiés en illisible),
- blocs oubliés (par ex. blocs de texte non segmentés),
- blocs parasites (par ex. blocs de texte sans contenu textuel),
- blocs composés (application des règles).



Les erreurs de segmentation détectées lors du contrôle automatique de structure (cf. section 7.1) ne sont pas comptabilisées ici.

Le taux d'erreur est calculé ainsi :

$$\text{taux d'erreur segm.} = \frac{\sum \text{erreurs segm.}}{\text{nombre de blocs}}$$

Il doit être inférieur au taux qualité attendu pour la segmentation (cf. CCTP).



Le mode opératoire de ce contrôle et le calcul du taux qualité seront détaillés dans une charte de contrôle OCR/ALTO élaborée conjointement par la BnF et le prestataire.

### 7.2.3 Détail des métriques qualité

La BnF et le prestataire s'accorderont sur le format descriptif à utiliser pour les indicateurs statistiques et les taux qualité demandés par la BnF.

*Dans le manifeste du document numérique (refNum ou METS)*

Les indicateurs à fournir (à l'échelle du document) sont notamment :

- le taux OCR brut : le taux de confiance en sortie de l'OCR sur la totalité de la page
- le taux OCR corrigé : le taux de confiance après correction manuelle et prise en compte des éventuels mots et zones déqualifiés,
- le coefficient d'étalonnage utilisé (cf. section 7.2.2)
- le taux NQA (niveau de qualité acceptable) moyen, exprimé par : *taux OCR corrigé × coefficient d'étalonnage*. C'est ce taux qui fait foi relativement aux exigences qualité du marché.

*Dans les fichiers ALTO*

Les indicateurs à fournir sont notamment, pour chaque page ALTO :

- nombre de blocs illisibles,
- nombre de mots illisibles,
- nombre de lignes dans les zones de courbures,
- nombre total de mots ALTO
- nombre total de caractères ALTO, par catégorie (lettres, chiffres, ponctuations)
- le taux OCR brut : le taux de confiance en sortie de l'OCR

- le taux OCR corrigé : le taux de confiance après correction manuelle (dernière ligne de l'exemple suivant). Ce taux sera identique au taux OCR brut dans le cas d'un document non corrigé ou d'une prestation sans correction.

Ce taux sera inséré dans l'attribut ACCURACY de l'élément Page (cf. section 5.1.3).

D'autres indicateurs pourront être mis en place, après étude avec le prestataire

- nombre de lignes en courbure détectées et traitées
- difficulté de post-correction de la page
- etc.

Ces indicateurs seront décrits à l'aide d'instances de l'élément <processing StepDescription>. Par exemple :

```
<OCRProcessing ID="OCR_1">
  <ocrProcessingStep>
    <processingDateTime>2013-04-23</processingDateTime>
    <processingStepDescription>[001_OCR_BRUT]OCR BRUT 96.240 %</processingStepDescription>
    <processingStepDescription>[002_CHARS]{12187}</processingStepDescription>
    <processingStepDescription>[003_CHARS_USED]{11262}</processingStepDescription>
    <processingStepDescription>[004_CHARS_SUSPECTS_USED]{0}</processingStepDescription>
    <processingStepDescription>[005_CHARS_UNUSED]{925}</processingStepDescription>
    <processingStepDescription>[006_CHARS_SUSPECTS_UNUSED]{143}</processingStepDescription>
    <processingStepDescription>[007_STRINGS]{2615}</processingStepDescription>
    <processingStepDescription>[011_CHARS_PUNCTUATION]{781}</processingStepDescription>
    <processingStepDescription>[012_CHARS_DIGITS]{52}</processingStepDescription>
    <processingStepDescription>[013_CHARS_SUSPECT]{143}</processingStepDescription>
    <processingStepDescription>[022_CHARS_UNUSED_BY_BLOCK]{0}</processingStepDescription>
    <processingStepDescription>[111_LIGS_DETECTED_CURV]{0}</processingStepDescription> ...
  </ocrProcessingStep>
</OCRProcessing>
```

Le tableau suivant présente la liste des indicateurs. Les indicateurs en gras doivent obligatoirement être inscrits dans les fichiers ALTO.

| Code                  | Description  | Exemple de valeur |
|-----------------------|--|-------------------|
| <b>001_OCR_BRUT</b>   | Taux OCR brut en sortie du moteur, calculé sur la totalité de la page  | OCR BRUT 96.240 % |
| <b>002_CHARS</b>      | Nombre total de caractères de la page  | 12187             |
| <b>003_CHARS_USED</b> | Nombre total de caractères utilisés pour le calcul du taux OCR pondéré (cf. section 4.1.3 : hors caractères en OCR brut, illisibles, chiffres, | 11262             |

|                             |   |      |
|-----------------------------|---|------|
|                             | ponctuations, spéciaux)   |      |
| 004_CHARS_SUSPECTS_USED     | Nombre total de caractères comptabilisés marqués comme suspects par l'OCR.<br>Toujours égal à 0   | 0    |
| 005_CHARS_UNUSED            | Nombre total de caractères non comptabilisés pour le calcul du taux OCR pondéré.<br>Rq : 005 = 002 - 003                                      | 925  |
| 006_CHARS_SUSPECTS_UNUSED   | Nombre total de caractères non comptabilisés marqués comme suspects par l'OCR<br>Rq : = tous les caractères des mots/blocs de type illisible  | 143  |
|                             | Rq : Nombre total de caractères en OCR brut = 002 - 003 - 006   | 782  |
| 007_STRINGS                 | Nombre total de mots (String)   | 2615 |
| 008_STRINGS_USED            | Nombre total de mots utilisés pour le calcul du taux OCR pondéré (hors mots en OCR brut et illisibles)  | 2544 |
| 009_STRINGS_SUSPECTS_UNUSED | Nombre total de mots non comptabilisés marqués comme suspects par l'OCR (illisibles isolés ou dans les blocs illisibles).                     | 21   |
|                             | Rq : Nombre total de mots en OCR brut = 007 - 008 - 009   | 50   |
| 010_CHARS_LETTERS           | Nombre de caractères de type lettre   |      |
| 011_CHARS_PUNCTUATION       | Nombre de caractères de ponctuation   | 781  |
| 012_CHARS_DIGITS            | Nombre de caractères de type chiffre arabe (0-9)  | 52   |
|                             | Rq : Autres caractères = 002 - 010 + 011 + 012  |      |
| 013_CHARS_SUSPECT           | Nombre de caractères suspects (comptabilisés ou non). Tous les caractères suspects sont ceux des textes typés en illisible.<br>Rq : 013 = 006 | 143  |
| 022_CHARS_UNUSED_BY_BLOCK   | Tous les caractères des blocs suspects (illisibles)<br>Rq : dans les blocs seuls, pas ceux des mots illisibles isolés                         | 100  |
|                             | Rq : Nombre de caractères suspects isolés = 013 - 022   | 43   |
| 023_BLOCKS_SUSPECTS         | Nombre total de blocs marqués comme suspects (illisibles) par   | 1    |

|                        |   |   |
|------------------------|---|---|
|                        | l'OCR.  |   |
| 111_LIGS_DETECTED_CURV | Nombre de lignes courbes détectées sur la page. | 0 |
|                        |   |   |
|                        |   |   |