



Europeana Newspapers

Contrôle de la qualité OCR

Christian Clausner, USAL

Traduction : Jean-Philippe Moreux, Bibliothèque nationale de France



PRImA

Pattern Recognition & Image Analysis
Research Lab



europeana
newspapers

Partie 1 : Introduction

Cas d'usage, workflow, outils, formats

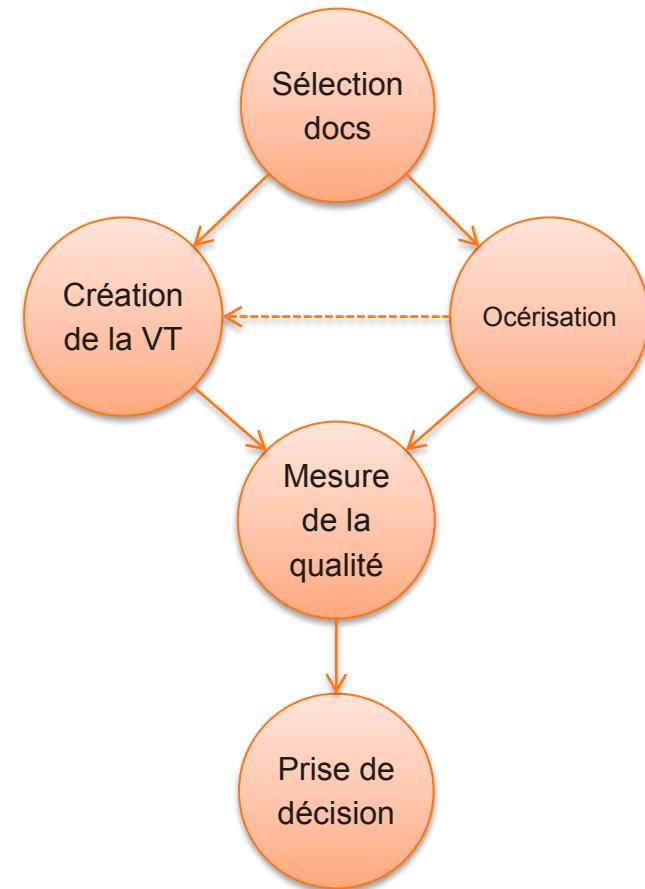
Contrôle de la qualité de l'OCR – Cas d'usage

- Etudes avant-projet
 - Évaluer si une collection d'imprimés est apte à être numérisée pour un usage particulier
- Contrôle de la qualité
 - Contrôler le résultat de l'OCR pendant ou après la production



Contrôle de la qualité de l'OCR – Workflow

1. Sélectionner un (petit) ensemble de documents représentatifs du corpus.
2. Océriser les documents dans la chaîne de numérisation cible.
3. Créer les documents de référence (vérité terrain, VT)
4. Mesurer la qualité en sortie de la chaîne de numérisation.
5. Prendre une décision en fonction des performances mesurées.

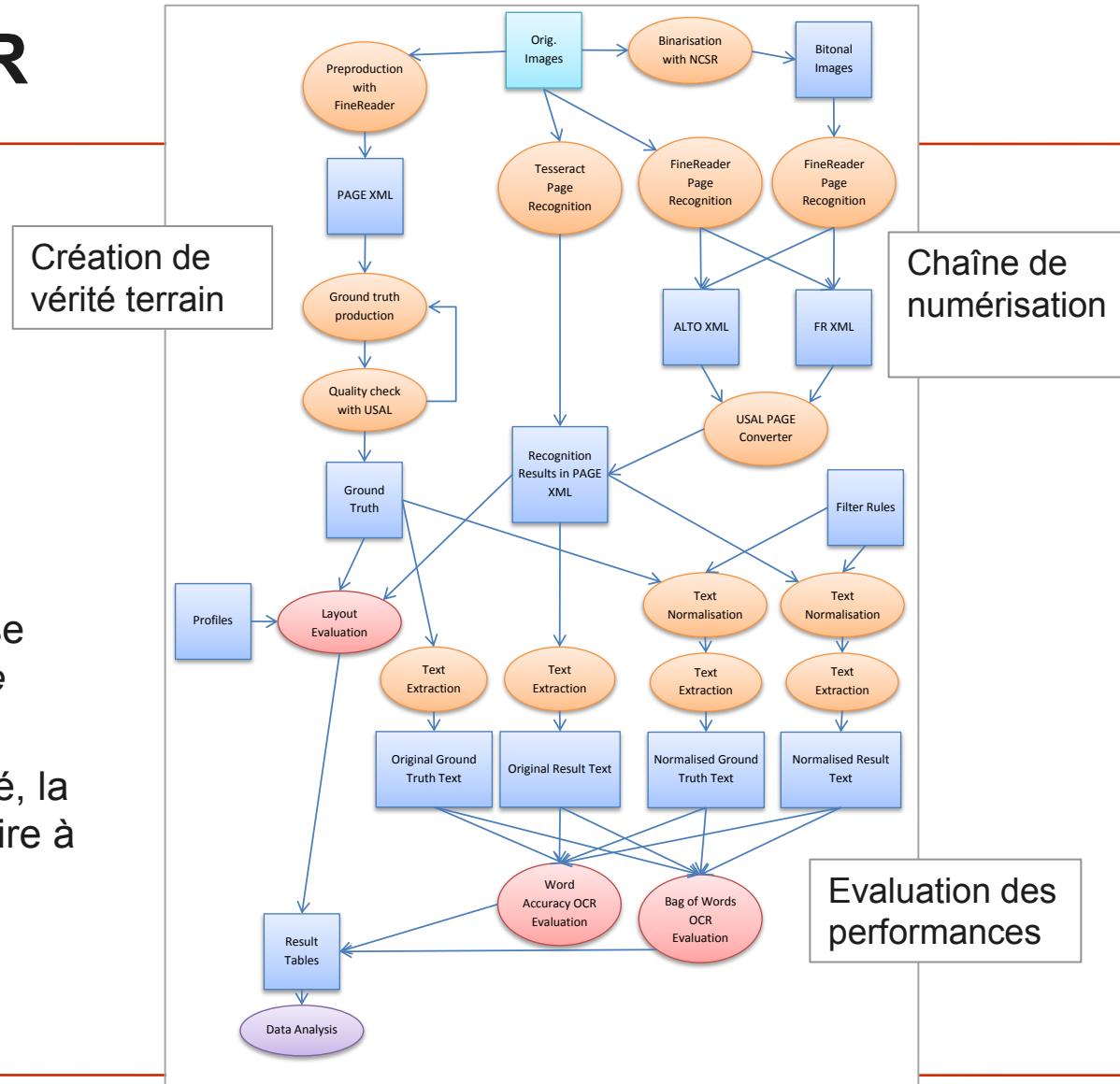


Contrôle de l'OCR

- Workflow

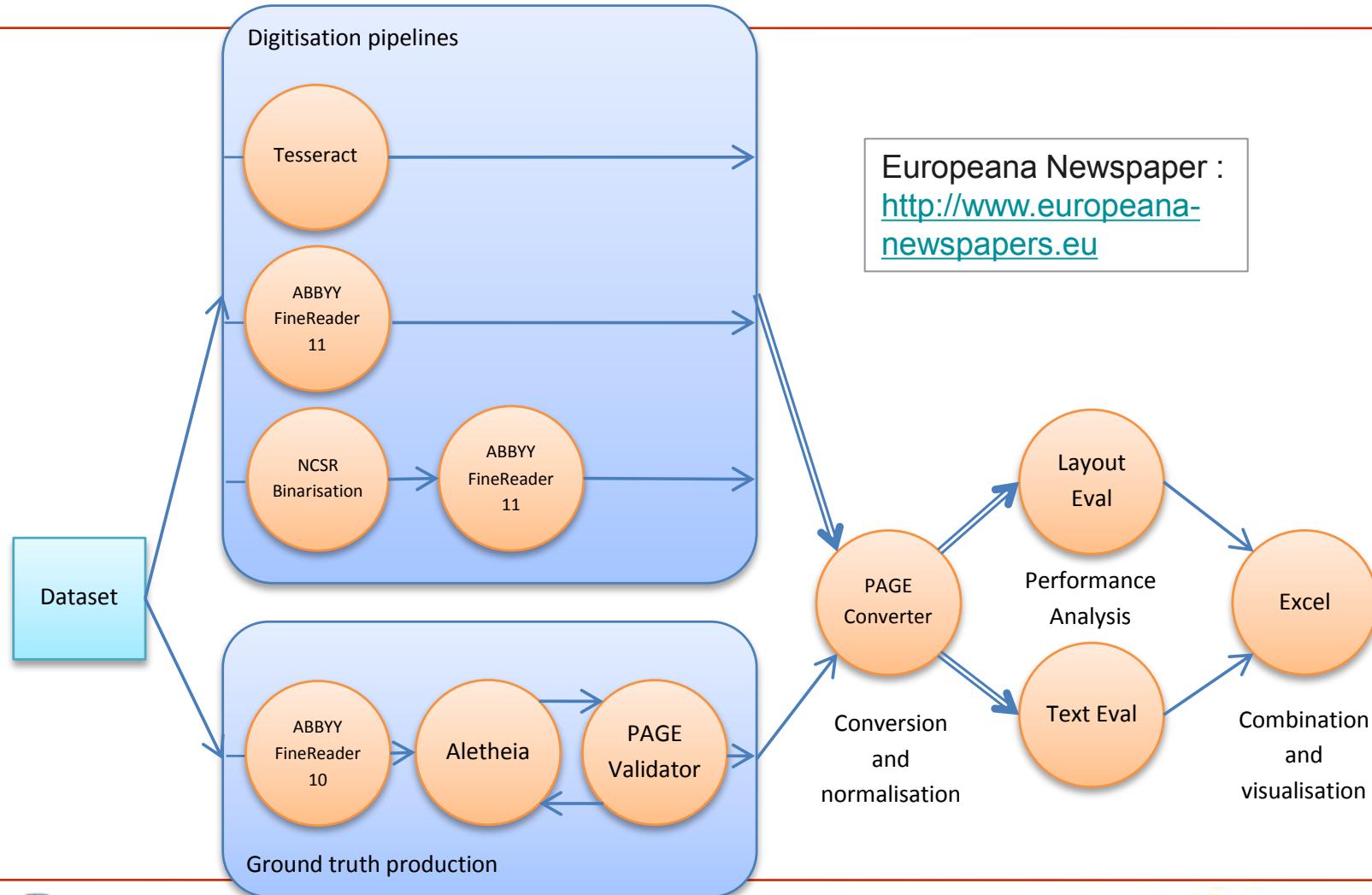
NB :

- Dans le cas d'une étude avant-projet, une chaîne de numérisation devra être mise en place (OCR open source ou autre).
- Pour un projet ou un marché, la chaîne est celle du prestataire à évaluer.



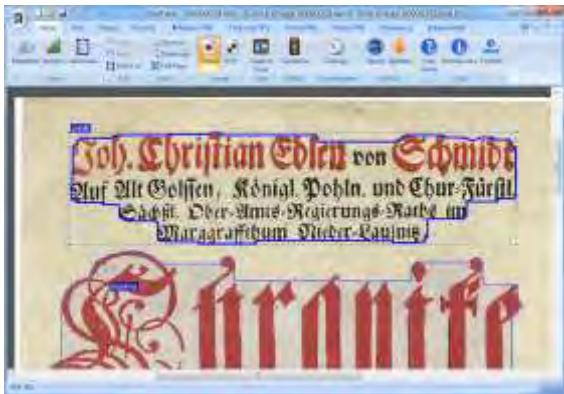
This project is partially funded under the ICT Policy Support Programme (ICT PSP) as part of the Competitiveness and Innovation Framework Programme by the European Community
http://ec.europa.eu/ict_psp

Exemple de workflow (Europeana Newspapers)

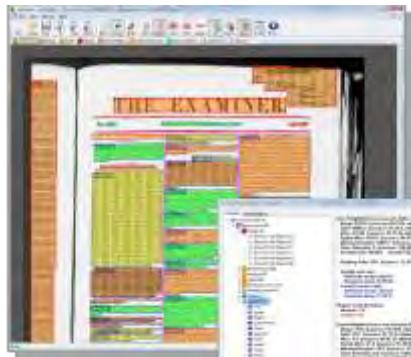


This project is partially funded under the ICT Policy Support Programme (ICT PSP) as part of the Competitiveness and Innovation Framework Programme by the European Community
http://ec.europa.eu/ict_psp

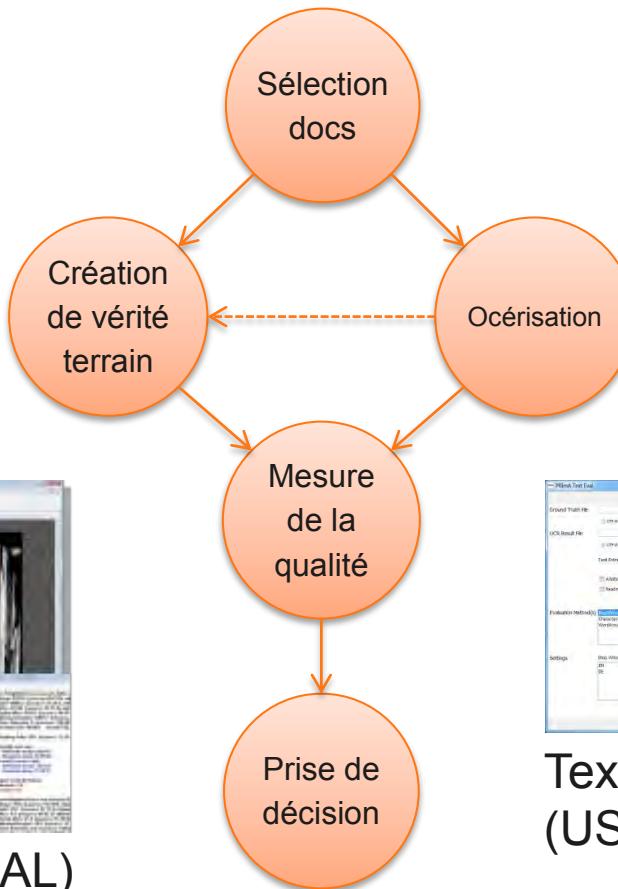
Contrôle de la qualité de l'OCR – Outils



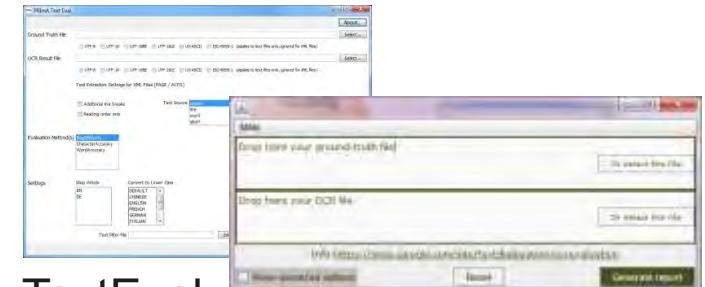
Aletheia



LayoutEval (USAL)



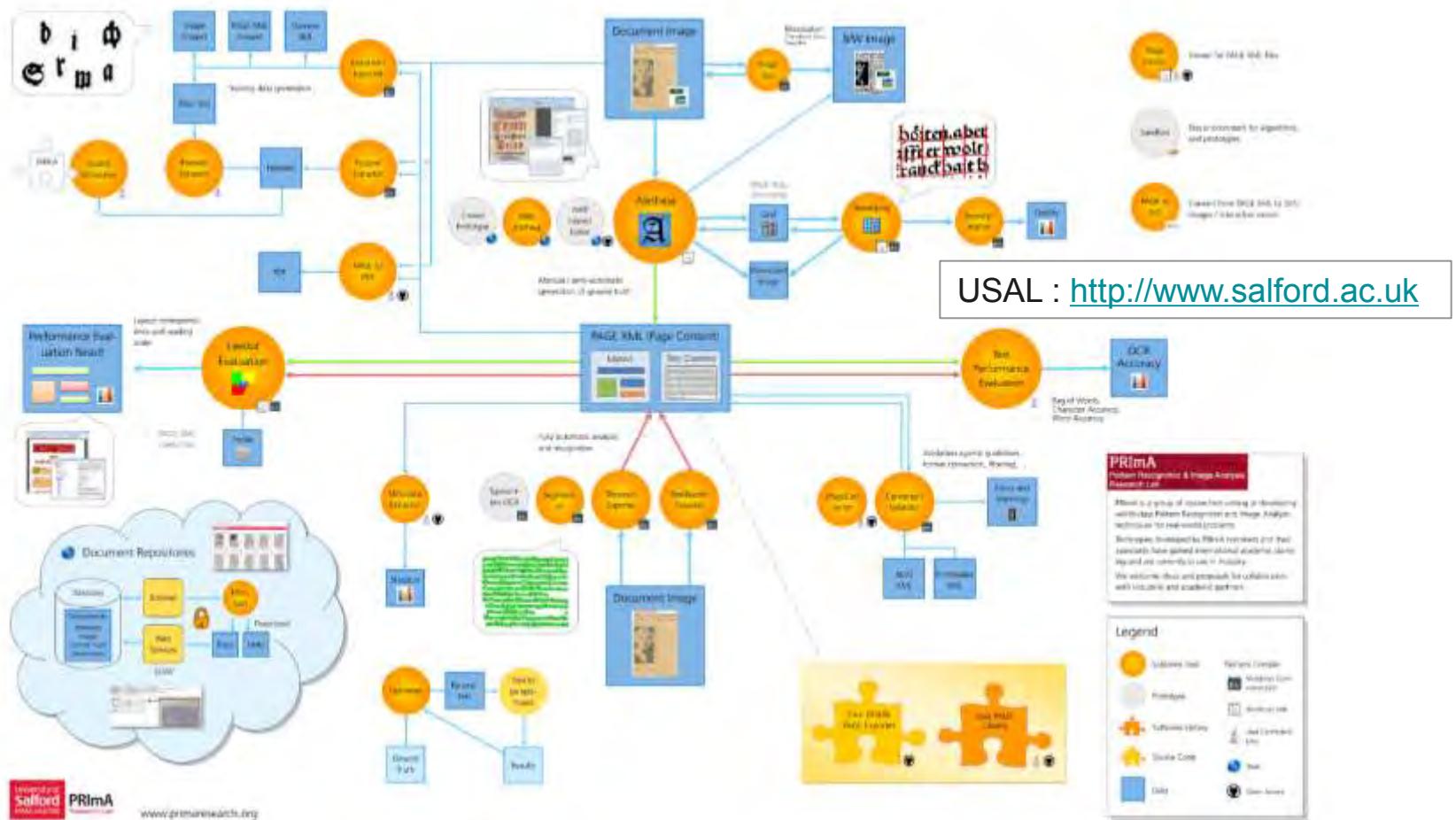
ABBYY FineReader
Tesseract
...



TextEval
(USAL)

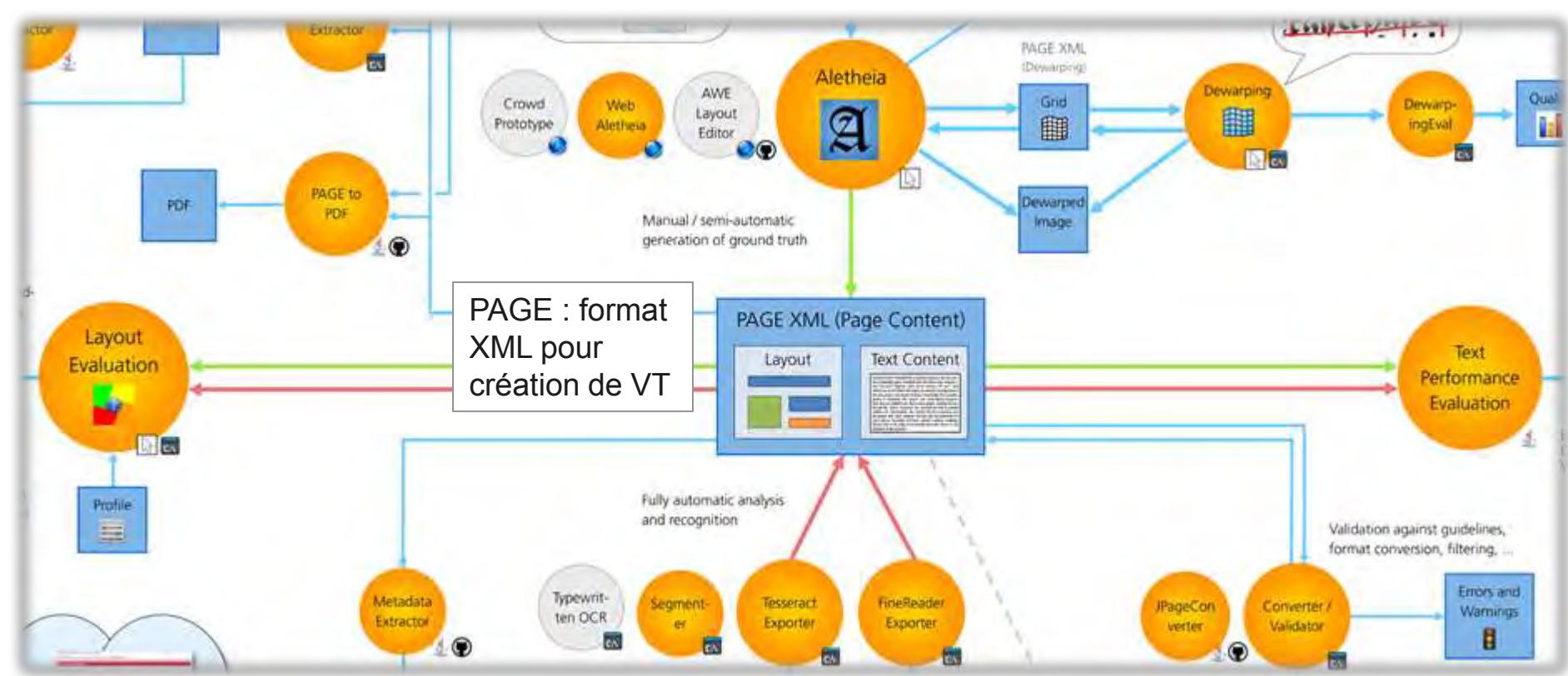
OCREvaluation
(Univ. Alicante)

Contrôle de la qualité de l'OCR – Outils et formats (USAL)



This project is partially funded under the ICT Policy Support Programme (ICT PSP) as part of the Competitiveness and Innovation Framework Programme by the European Community
http://ec.europa.eu/ict_psp

Contrôle de la qualité de l'OCR – Format PAGE



This project is partially funded under the ICT Policy Support Programme (ICT PSP) as part of the Competitiveness and Innovation Framework Programme by the European Community
http://ec.europa.eu/ict_psp



europeana
newspapers

Partie 2 : Océrisation

Techniques et outils

OCR – Analyse de page et reconnaissance

Image



Numérisation

- Segmentation
- Classification des natures de contenu
- OCR
- ...

OCR



ABSENT YET PRESENT. by Gilberta M. F Lyon, 3 vols., London; Digby, Long, and Co. We were able to speak in very high praise of Miss Lyon's former story, "For Good or Evil." This, however, is far better and must be classed amongst the best books of the season. The writer has conspicuous talents; in her description of the scenery whither she guides us in the portraiture of the characters with which she peoples her story, in the varied pleasantries and troubous periods which go to make up the delightful romance which she unfolds in her pages, we are enchain'd by the most sympathetic ties, until the fortunes of hero and heroine are spun out to their

OCR – Alternatives

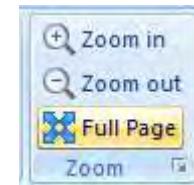


- ABBYY FineReader Engine (FRE)
- Tesseract (*open source*)
- Aletheia (avec Tesseract intégré)



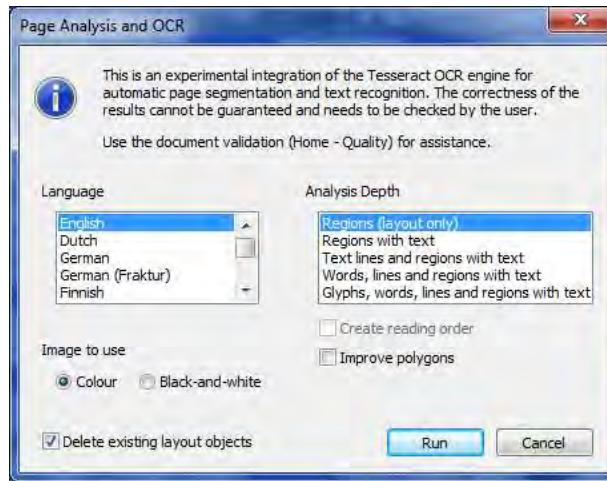
OCR – En pratique avec Aletheia

1. Créer un nouveau document.
2. Sélectionner **IMG.tif**
3. Confirmer “... without B/W image”.
4. Zoom “Full page”.
5. Pour lancer l’OCR : “Analyse Page”.



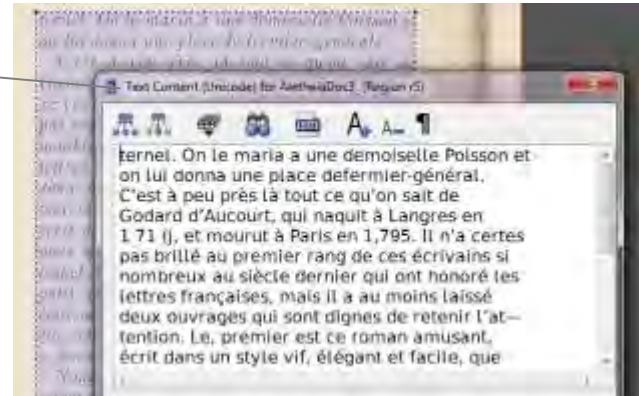
OCR – En pratique avec Aletheia

1. Paramétriser.
2. Cliquer sur “Run”.
3. Attendre...
4. Visualiser.

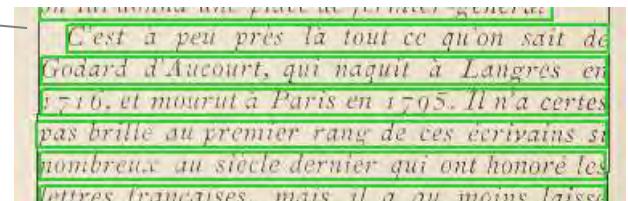
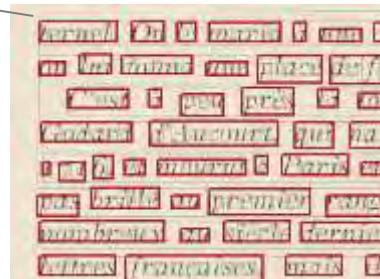


OCR – En pratique avec Aletheia

1. Consulter le texte océrisé : F11 ou activer “Texte overlay”



2. Visualiser ligne de texte (F7), mots (F8) et glyphes (F9)



OCR – En pratique avec Aletheia

1. Enregistrer le document au format PAGE.



2. Si besoin, exporter le texte.



europeana
newspapers

Partie 3 : Création de vérité terrain

Cas d'usage, méthodes, outils

L'OCR est-il bon ?

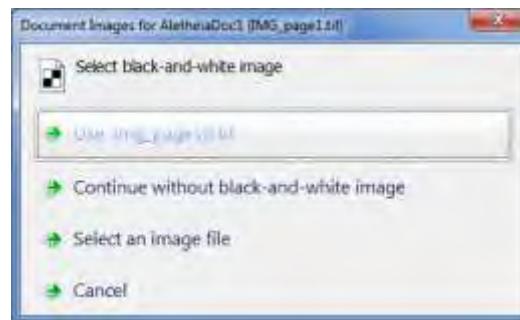


- Un document a été océrisé.
- Mais est-il de bonne qualité ?
- Un contrôle visuel est possible mais subjectif...
- Solution : comparer avec une référence

→ Vérité terrain

Création de vérité terrain – En pratique avec Aletheia

- Aletheia permet de créer un OCR de référence à partir d'une image couleur
- et éventuellement d'une image déjà binarisée :



- NB : il est aussi possible de créer l'image binarisée dans Aletheia.

Aletheia :
<http://www.primaresearch.org/tools>

Aletheia – Outils niveau image

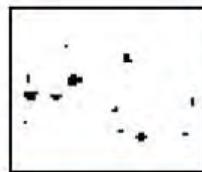


Chargement d'images (TIFF, PNG, JPEG)

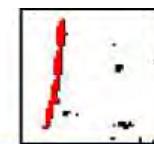
Binarisation d'image :

- seuillage
- Otsu (automatique)
- Sauvola (automatique, adaptative)

Suppression du bruit



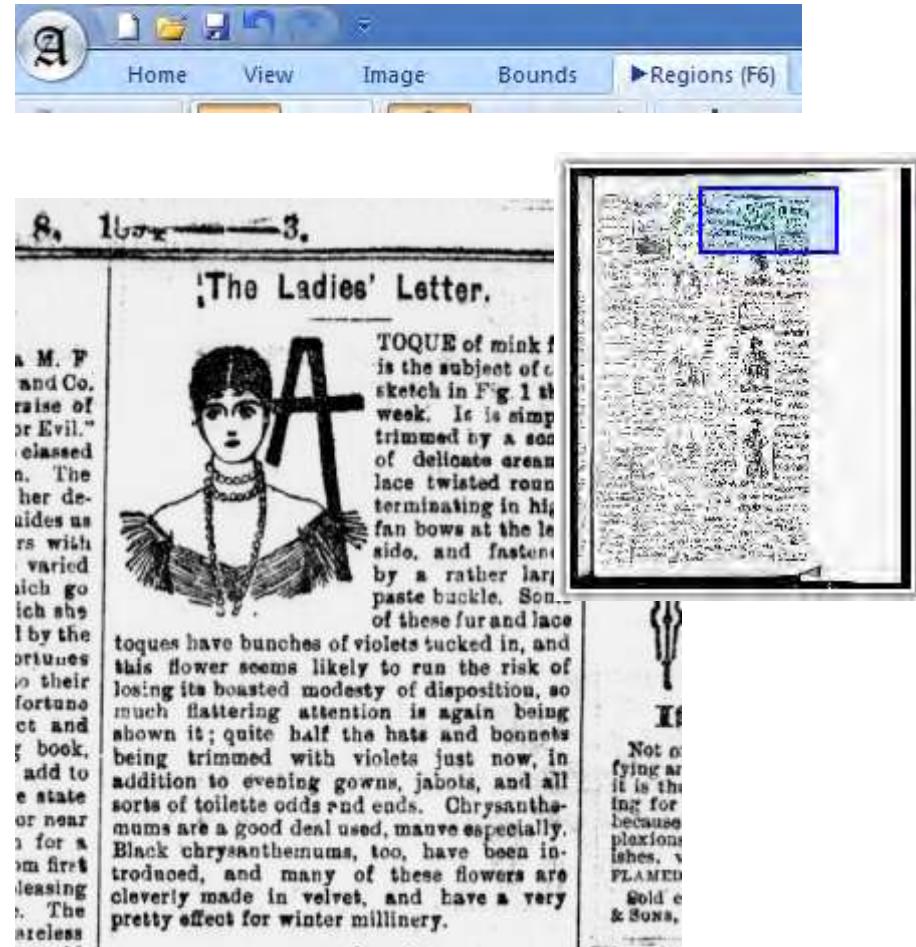
Suppression manuelle du bruit



Suppression des bordures (sur image N&B)

Aletheia – Création manuelle de régions

1. Onglet “Regions” (F6)
2. Zoomer sur une zone.
3. Outil “Rectangle”.
4. Tracer une boîte (ou un polygone) autour d'un bloc de texte.



This project is partially funded under the ICT Policy Support Programme (ICT PSP) as part of the Competitiveness and Innovation Framework Programme by the European Community
http://ec.europa.eu/ict_psp

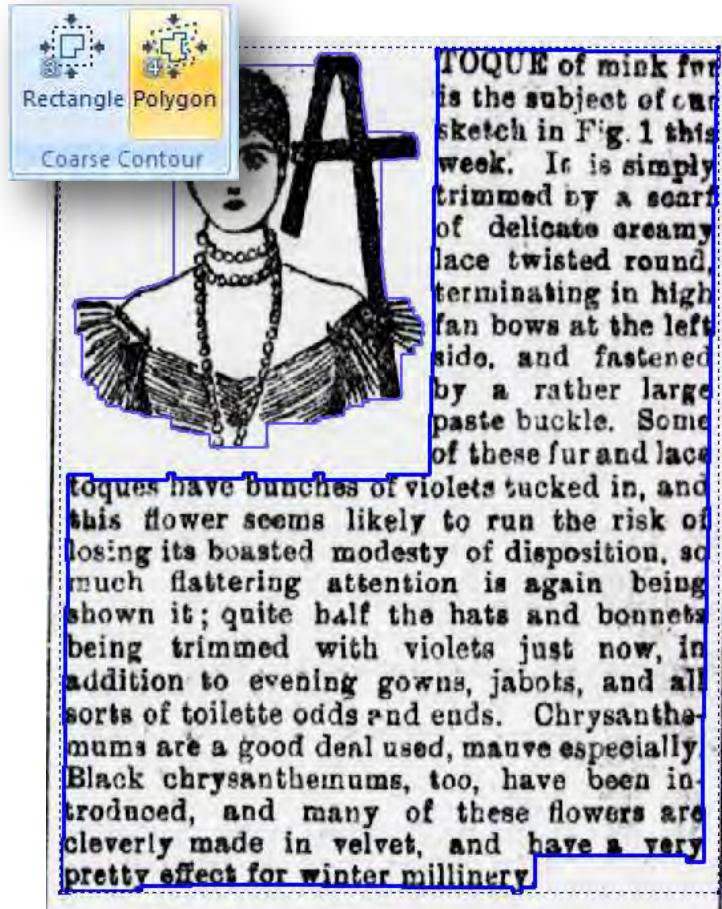
Aletheia – Cr ation semi-automatique de r gions

1. Outil “Fine Contour” -
“Rectangle” (*contour fin*)
2. Tracer une bo te autour de
l’image (tout inclure, mais
pas au-del )
3. Annuler et
essayer de nouveau
si n cessaire.



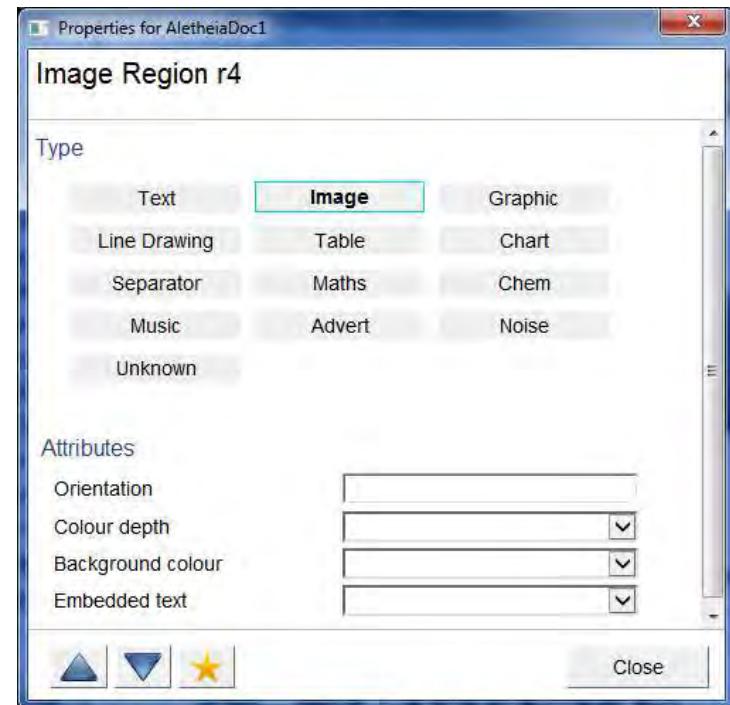
Aletheia – Cr eation semi-automatique de r egions

1. Outil “Coarse Contour” -
“Polygon” (*contour grossier*)
2. Tracer une bo te autour d’un
paragraphe de texte (un clic
ajoute un point ; clic droit
pour finir)



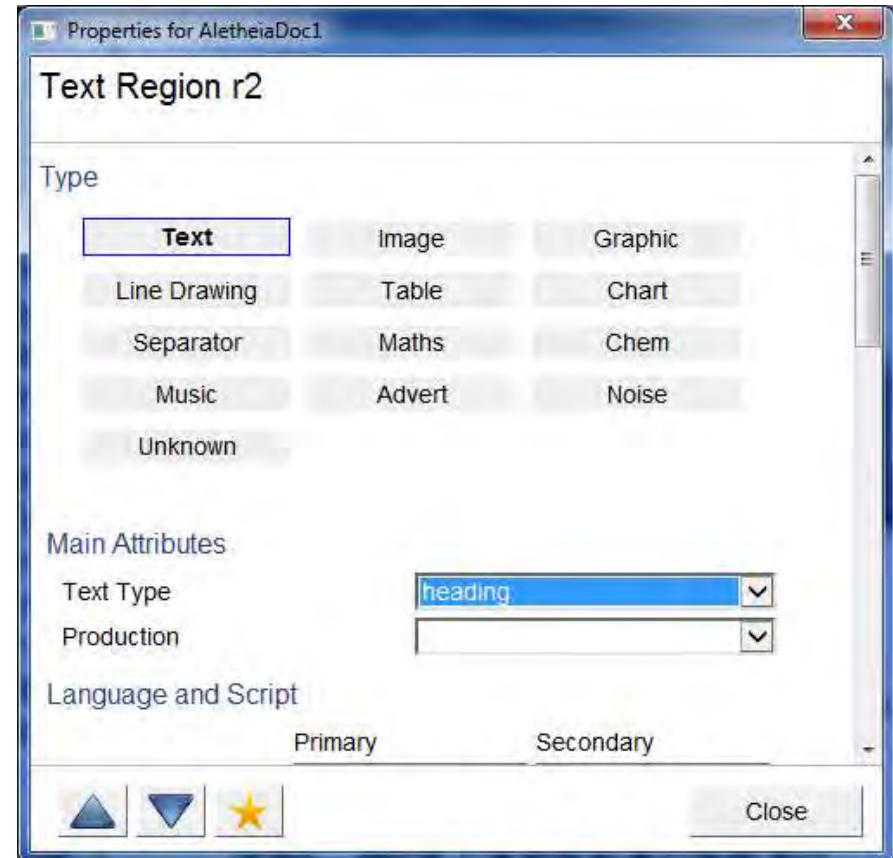
Aletheia – Description des régions / illustration

1. Double-clic sur une région illustration (ou F10).
2. Choisir “Image” dans la boîte de dialogue pour changer le type de région.
3. Renseigner les attributs en fonction du type de région.



Aletheia – Description des régions / texte

1. Cliquer sur un titre. F10.
2. Changer “Text Type” de “Paragraph” en “Heading”
3. Les icônes Haut et Bas permettent de naviguer de bloc en bloc.
4. Fermer le dialogue.



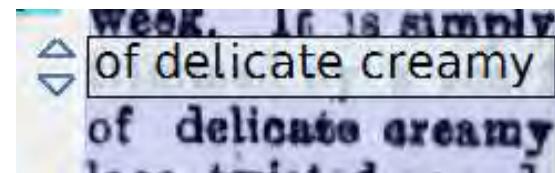
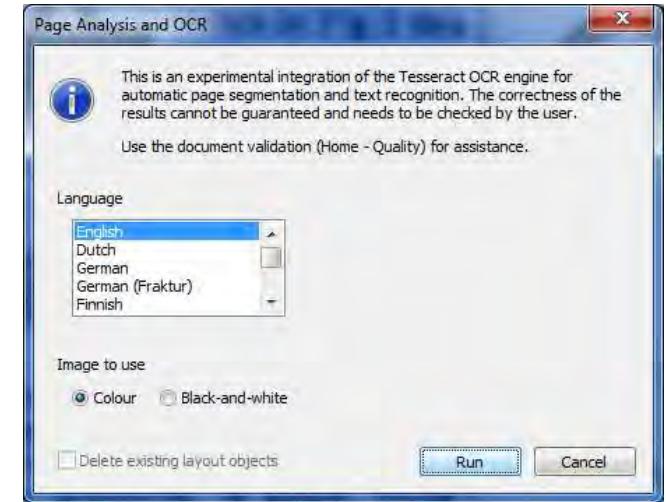
Aletheia – Saisie du texte

1. Cliquer sur un titre.
2. Ouvrir la boîte de dialogue “Text” (F11)
3. Saisir le texte du bloc.
4. Des caractères spéciaux sont disponibles via le clavier virtuel.
5. Cliquer sur “Next” pour passer à la saisie du bloc suivant.



Aletheia – Creation de regions – OCR

1. Cliquer sur “OCR Region”.
2. Choisir l’image couleur ou N&B.
3. En fonction de la qualite, il est plus rapide de corriger ou de ressaisir...
4. L’option “Text overlay” permet de visualiser le texte au survol de souris.
5. Les touches Haut/Bas permettent un controle ligne a ligne.



Aletheia – Ordre de lecture

L'ordre de lecture logique doit être défini dans la VT. Il est important pour certains usages du document numérique : recherche de phrases complète, lecture écran ou en synthèse sonore, conversion en livre numérique, etc.

Aletheia permet de créer l'ordre de lecture manuellement ou par analyse de la page.



Aletheia – Cration manuelle de l'ordre de lecture

1. Ouvrir la fentre “Reading order and Layers” (F12)
2. Slectionner “Group (ordered)” – C’est la racine du premier groupe de contenu du document.
3. Double-clic sur le premier bloc de contenu.
4. Double-clic sur le bloc suivant pour crer l’ordre de lecture. Etc.

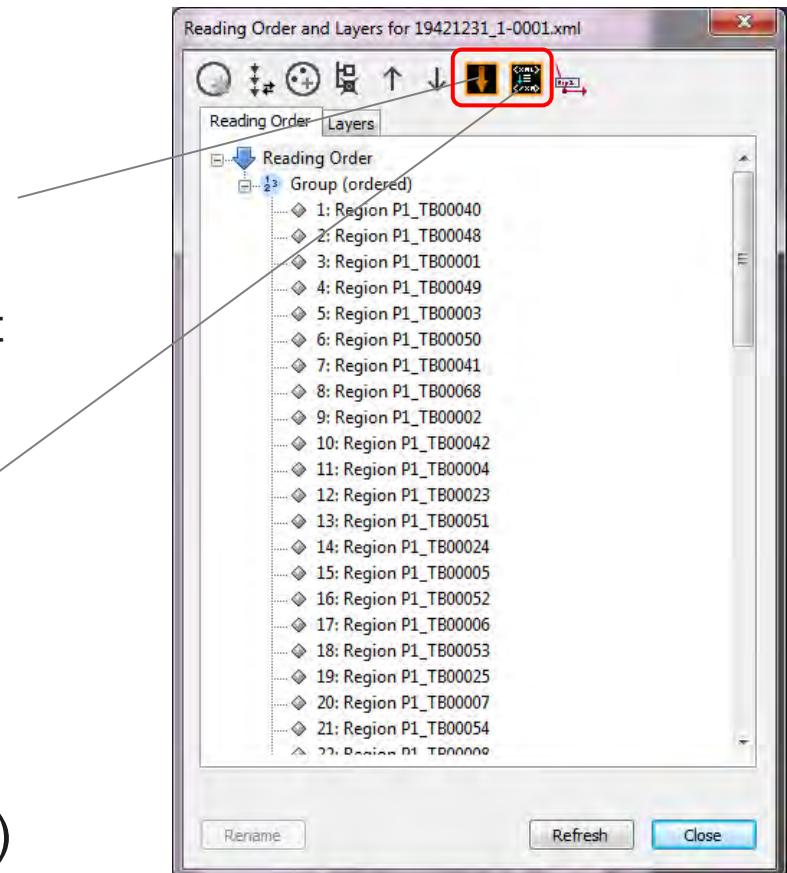
Alternative : placer tous les contenus dans un groupe ou crer des groupes et de dfinir leur ordre de lecture.



Aletheia – Cration automatique de l'ordre de lecture

1. Dans la fentre “Reading order and Layers” (F12), deux options :

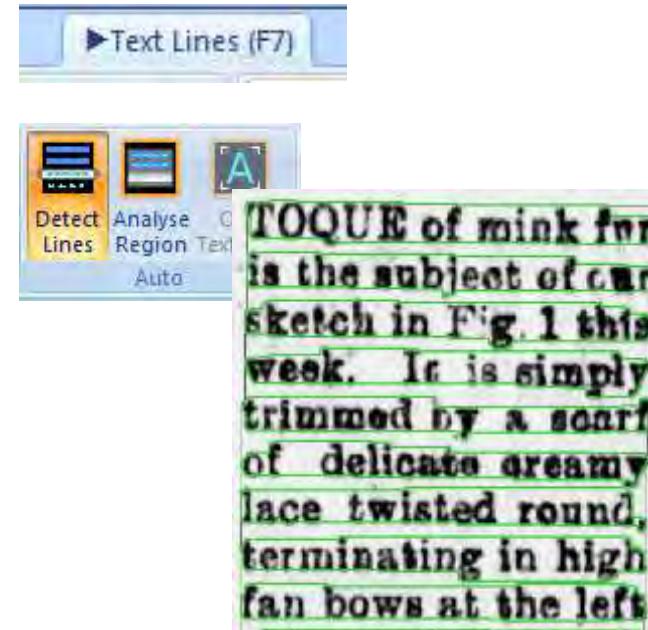
- “Create top-to-bottom reading order” – Analyse descendante et gauche-droite des blocs de la page. NB : Ne convient pas pour le multicolonnage.
- “Create reading order from internal region list (as in XML)” – Si l’ordre de lecture est djà prsent dans le fichier (cas de la cration de VT  partir d’un OCR, voir page 34).



2. L’ordre peut tre corrig (flches )

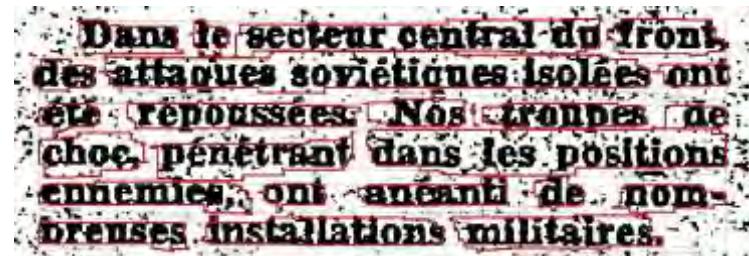
Aletheia – Cration des lignes de texte

1. Sélectionner l'onglet “Text Lines” (F7).
2. Clic sur “Detect Lines”.
3. Clic sur un paragraphe de texte. Les lignes sont détectées et créées.
4. Pour retoucher les lignes : outils manuels ou automatiques



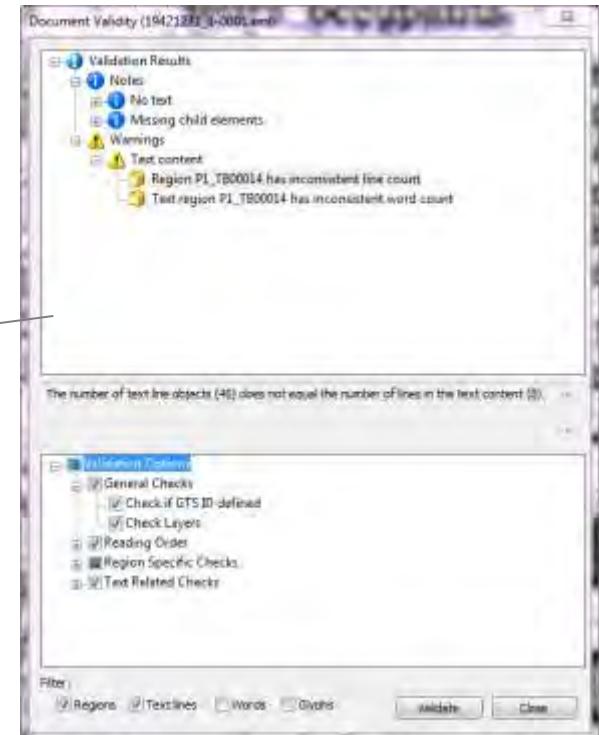
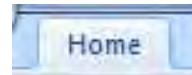
Aletheia – Mot, caractère

1. Sélectionner l'onglet “Words” (F8).
2. Clic sur “Detect Words”.
3. Clic sur un paragraphe de texte. Les mots sont détectés et créés. La segmentation peut alors être corrigée.
4. Même principe avec les caractères (onglet “Glyphs”, F9).



Aletheia – Validation

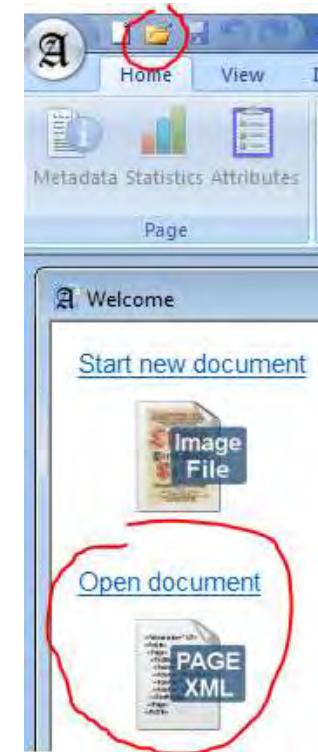
1. Sélectionner l'onglet “Home”.
2. Clic sur “Validation”.
3. Clic sur “Validate”. Il est possible de cibler les contrôles (bloc, ligne, mot, caractère)
4. Les résultats de validation sont affichés :
 - Information
 - Avertissement



Aletheia – Utiliser un OCR

La création ex nihilo d'une vérité terrain étant laborieuse, il est possible d'utiliser un document OCR comme base. Aletheia lit les formats Abbyy FineReader XML et ALTO XML.

1. Pour ouvrir un fichier OCR : “Open document”
2. Sélectionner l'image associée.
3. Enregistrer le document au format PAGE.

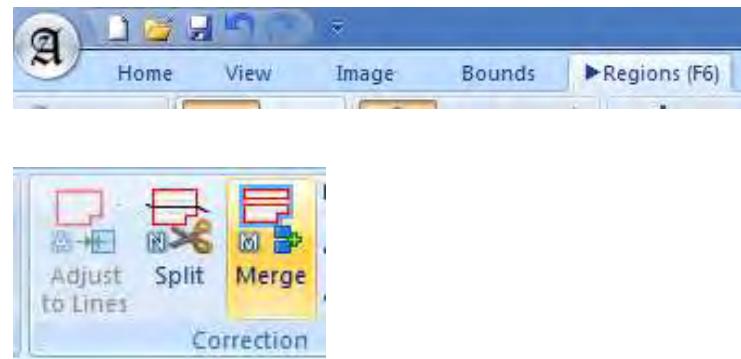


Aletheia – Utiliser un OCR

La segmentation des blocs peut ensuite être corrigée, ainsi que le texte.

- Edition des contours de bloc, découpage et fusion de blocs.
- Correction du texte.
- Correction de l'ordre de lecture.

NB : selon l'usage attendu pour la vérité terrain, la granularité de la segmentation peut varier. Par exemple pour le contrôle de la transcription du texte, le niveau bloc suffit.



Aletheia – Conversion de formats : OCR vers PAGE



ABBYY FineReader Engine (FRE)

- FR XML ▶ PageConverter (USAL)
- ALTO ▶ PageConverter (USAL)
- FineReader Integration



Tesseract (*open source*)

- TesseractToPage (USAL)

PAGE XML



europeana
newspapers

Partie 4 : Évaluation de la reconnaissance du texte

Cas d'usage, méthodes, outils

Evaluation des performances

Vérité
terrain



OCR
à évaluer



Mesure
de la
qualité



Evaluation de la reconnaissance du texte – Objectif

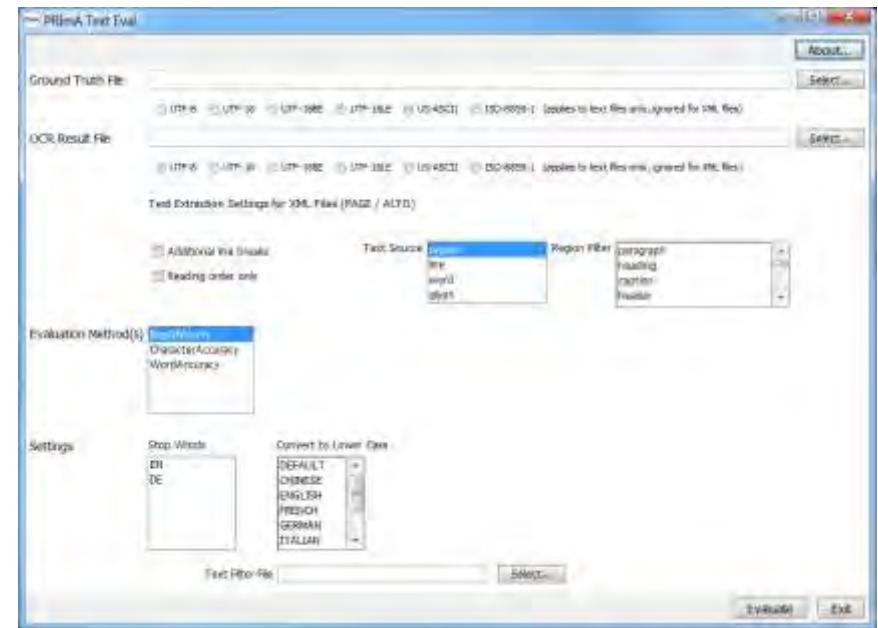
- Comparer le texte de la page avec le texte de référence.
- Différentes mesures classiques :
 - Précision niveau “caractère” (nombre des caractères erronés)
 - Précision niveau “mot” (nombre des mots erronés)
 - “Sac de mots” : l’ordre de lecture ne compte pas

Note : normalisation du texte

- Pour une évaluation moins stricte, remplacer certains caractères par leur version “normalisée”, par ex. :
 - Accent é → e
 - Ligature æ → ae
 - Tiret long – → -
 - S long (ſ) → s

Reconnaissance du texte – En pratique/Text Eval

- Outil **Text Eval** (USAL)
- Fichiers (texte ou XML) :
 - vérité terrain (“Ground truth”)
 - OCR à évaluer (“OCR result”)
- Choix de la méthode d'évaluation
- Options :
 - filtre d'éléments XML
 - liste de mots vides
 - ignorer la casse des caractères



Text Eval :

<http://www.primaresearch.org/tools/PerformanceEvaluation>

Reconnaissance du texte – En pratique/Text Eval

Evaluation en mode « Sac de mots »



Evaluation Method	Bag of Words	Word Accuracy
Measure		Value
wordIndexMissErrorRate		0.3500432152117545
wordIndexFalseDetectionErrorRate		0.47191011235955055
wordIndexSuccessRate		0.5827198760170476
numberOfUniqueWordsInGroundTruth		2314
numberOfUniqueWordsInResult		2848
wordCountMissErrorRate		0.3979342580614643
wordCountFalseDetectionErrorRate		0.4399146717231571
WordCountSuccessRate		0.5803173053145178
wordCountPrecision		0.5600853282768429
wordCountRecall		0.6249669399629728
wordCountFMeasure		0.5907500000000001
numberOfWordsInGroundTruth		3781
numberOfWordsInResult		4219

taux qualité
= 58 %

Reconnaissance du texte – En pratique/Text Eval

Evaluation « mots »

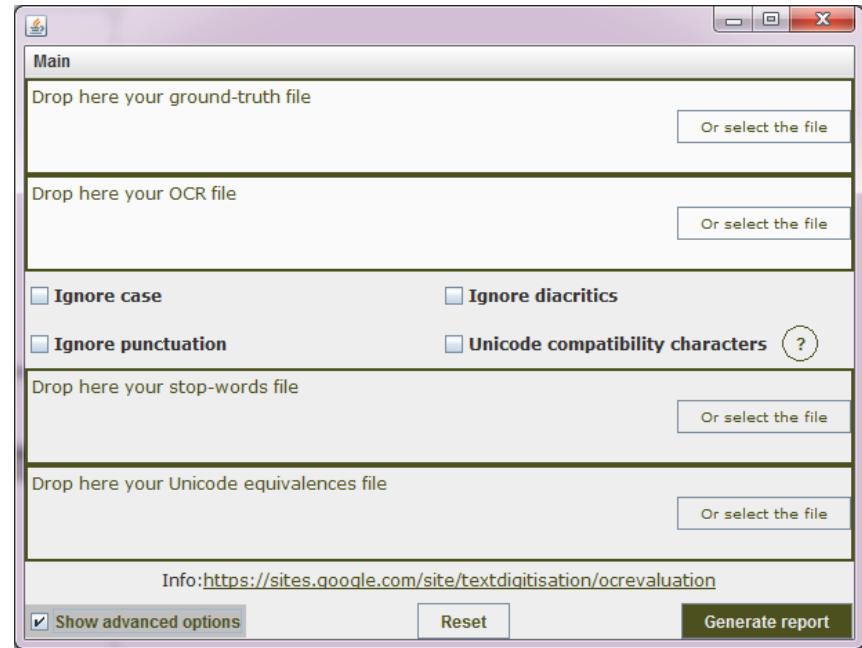


Evaluation Method	Bag of Words	Word Accuracy
Measure		Value
wordsInGroundTruthExclStopWords	3781	
wordsInResultExclStopWords	4219	
wordAccuracyExclStopWords	0.29780481354139116	
wordsInGroundTruth	5980	
wordsInResult	6195	
wordAccuracy	0.4508361204013378	

taux qualité
= 45,1 %

Reconnaissance du texte – En pratique/ocrevalUAtion

- **Outil ocrevalUAtion**
(université d'Alicante)
- Fichiers (texte ou XML)
- Options :
 - filtre d'éléments XML
 - liste de mots vides
 - équivalence de caractères
 - ignorer la casse et l'accentuation des caractères, les caractères de ponctuation



OCREvaluation:
<https://github.com/impactcentre/ocrevalUAtion>

Transcription du texte – En pratique/ocrevalUAtion

CER : taux erreur caractères (2,75 %)
WER : taux erreur mots (13,27 %)
WER : taux erreur « sac de mots » (12,39 %)

General results

CER	2,75
WER	13,27
WER (order n-gramme)	12,39

Difference spotting

D0000014.txt	X0000014.xml
Pour le mariage de la Religion , cette charité pour les pauvres , c'est évidemment à tous les devoirs du Saint Ministère & toutes ces vertus paternelles qu'vous caractérissez Si vous résidiez au Royaume de nos hommes , elles nous font domer . MONSIEURHEU , de jour long temps des gouvernements sont dans , que les Habitans de Bayeux , l'église primaient dans la cathédrale leur est la plus précieuse , j'aurai d'espérer . MONSIEURHEU , que l'Histoire de notre Ville vous sera agréable , Père & Fille -tous de Chavres qui vous ont aussi chez , vous êtes avec place les personnes qui y sont ar- rivés . Je suis avec un très profond respect . MONSIEURHEU , DE VOTRE GRANDEUR . La très-humble & très- obéissant Serviteur BELLER	

Error rate per character and type

Character	Char code	Total	Spurious	Corrigible	Lost	Error rate
À	128	1	0	1	0	1,56
À	26	4	0	4	0	0,00
À	27	4	0	4	0	0,00
À	28	13	0	13	0	53,85
À	29	5	0	5	0	0,00
À	29	2	0	2	0	50,00
À	30	1	0	0	0	0,00
À	41	1	0	1	0	0,00
À	42	2	0	2	0	0,00
À	43	1	0	1	0	0,00
À	44	2	0	2	0	0,00
À	45	11	0	11	0	0,00
À	47	2	0	2	0	0,00
À	48	2	0	2	0	0,00
À	49	8	0	8	0	0,00
À	49a	1	0	1	0	0,00
À	50	1	0	1	0	0,00

Alignement
des deux textes



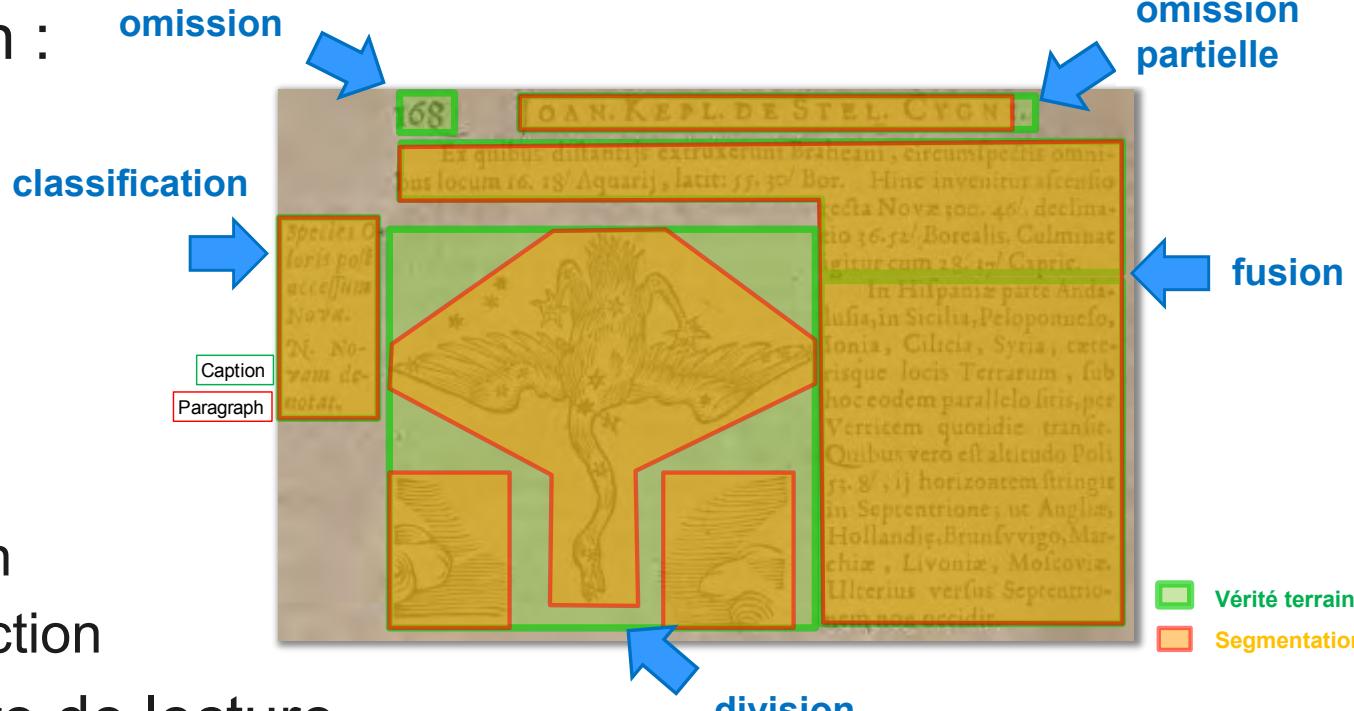
europeana
newspapers

Partie 5 : Évaluation de la segmentation

Cas d'usage, méthodes, outils

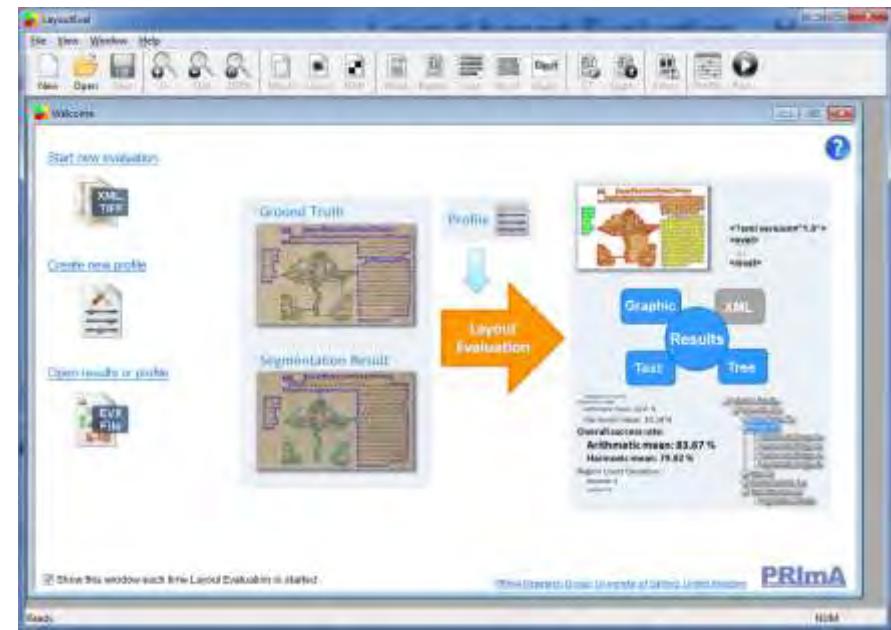
Évaluation de la segmentation

- Six types d'erreur de segmentation : omission
 - omission
 - omission partielle
 - division
 - fusion
 - défaut de classification
 - fausse détection
- Défaut d'ordre de lecture



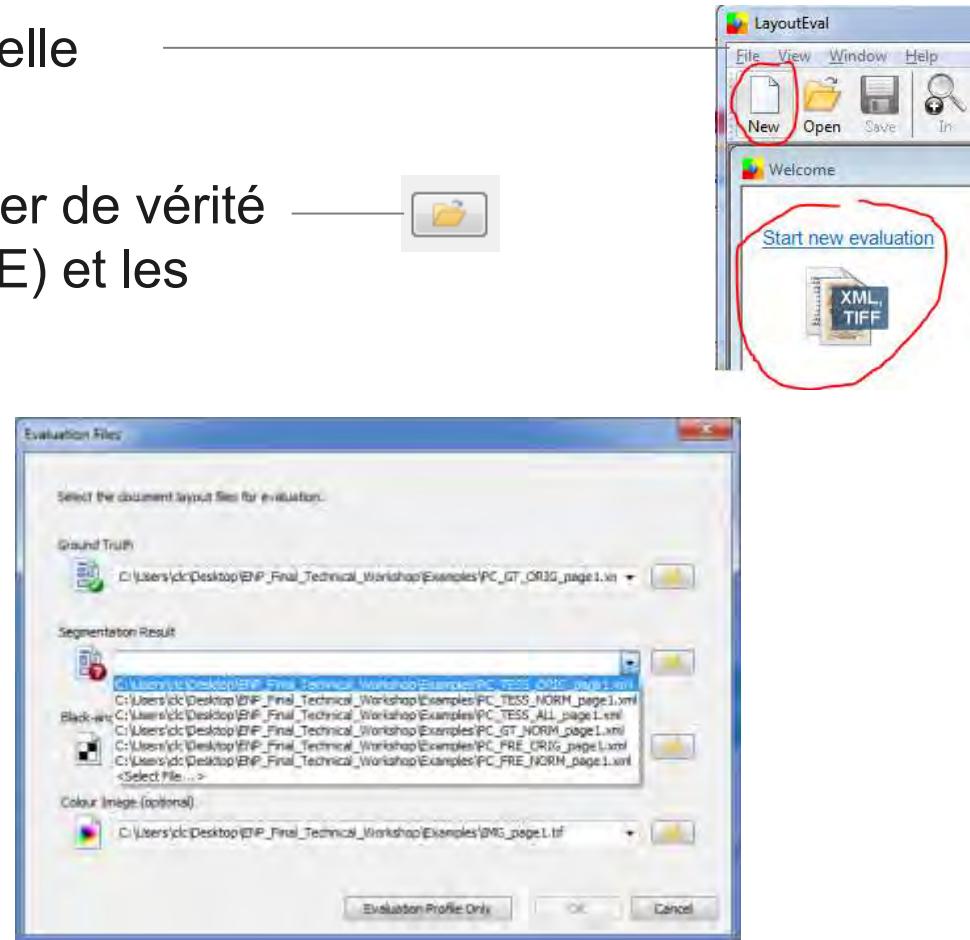
Évaluation de la segmentation – En pratique

- Outil : LayoutEval (USAL)
- Fichiers :
 - Vérité terrain (format PAGE)
 - Résultat de l'OCR (ALTO, Fine Reader, Tesseract)
 - Image binarisée
 - Image couleur (optionnelle)



Évaluation de la segmentation – En pratique

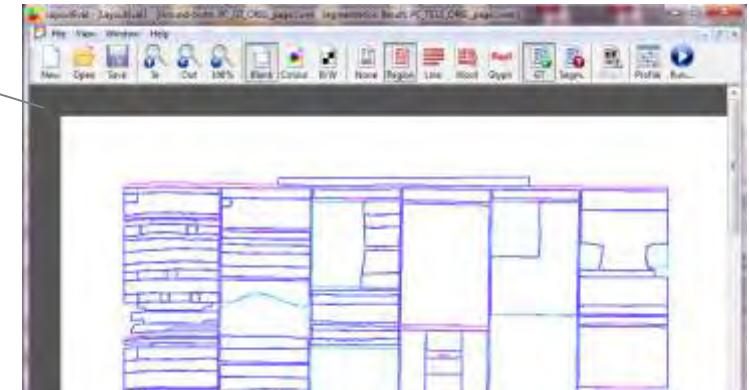
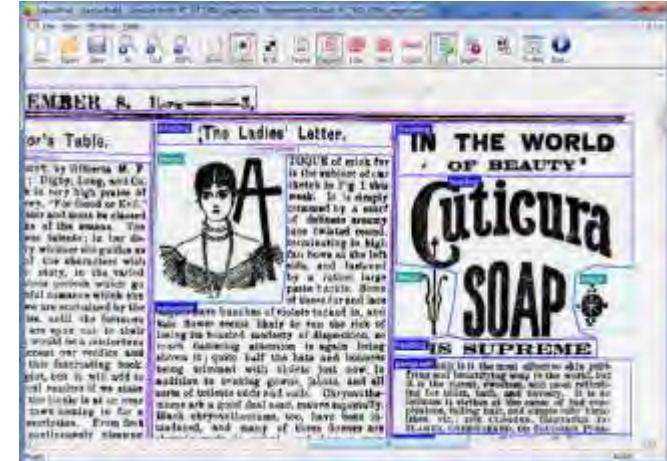
1. Démarrer une nouvelle évaluation.
2. Sélectionner le fichier de vérité terrain (format PAGE) et les images.
3. Sélectionner le fichier OCR.
4. Clic sur “OK”



This project is partially funded under the ICT Policy Support Programme (ICT PSP) as part of the Competitiveness and Innovation Framework Programme by the European Community
http://ec.europa.eu/ict_psp

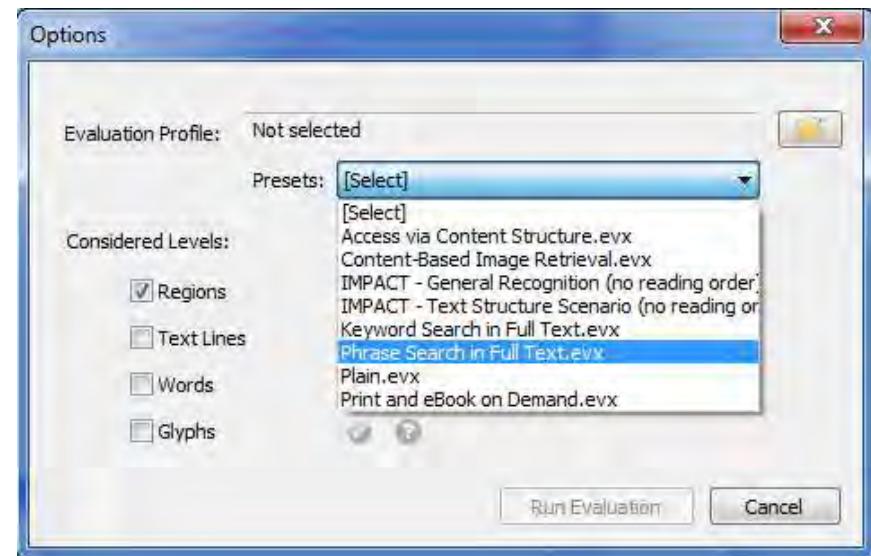
Évaluation de la segmentation – En pratique

1. Activer la vue “Region”
2. Se déplacer dans la page
3. Clic sur les boutons “GT” (vérité terrain) et Segm. pour visualiser vérité terrain, segmentation, ou les deux
4. Blank, Color et B/W : changer le mode de visualisation (Blank = segmentation seule)



Évaluation de la segmentation – En pratique

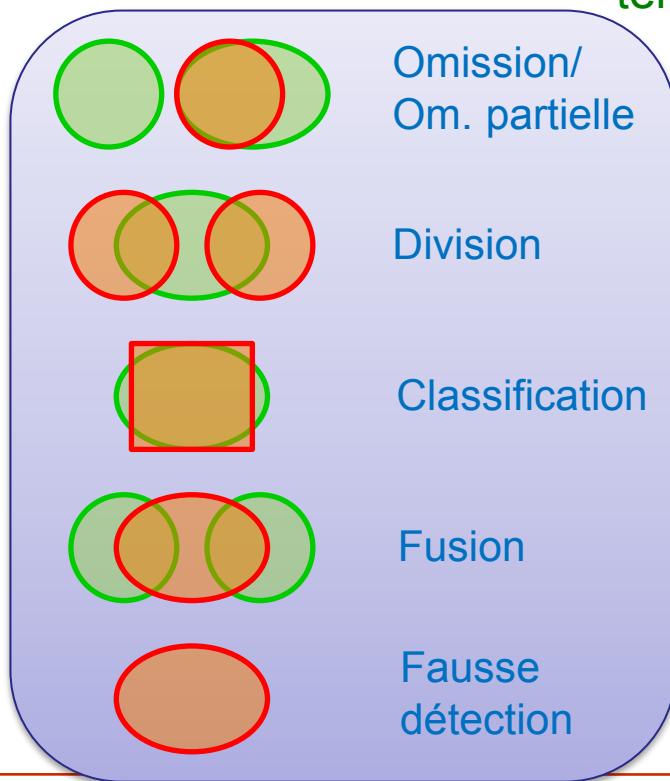
1. Clic sur “Run...” 
2. Sélectionner le profil d'évaluation prédefini "Phrase search in full text" et le niveau "Regions"
3. Clic sur "Run Evaluation"
4. Attendre...



While you are waiting...

Les métriques de l'évaluation sont basées sur l'analyse du recouvrement des régions

Types d'erreurs

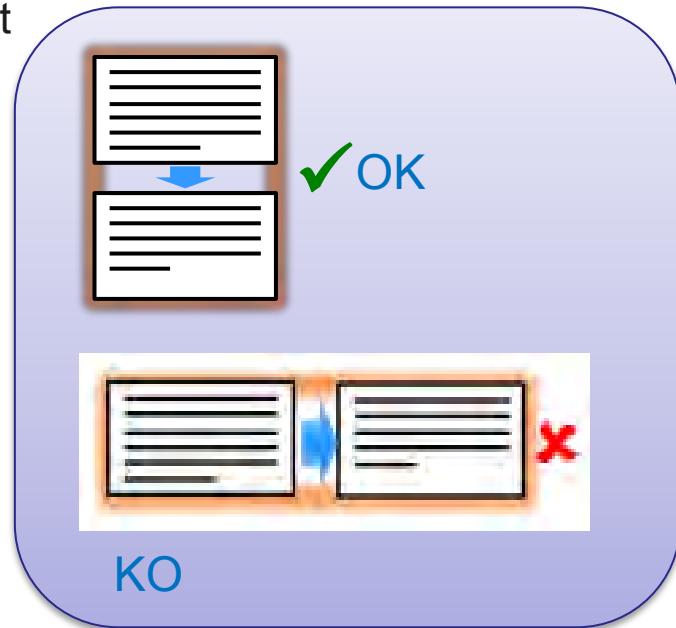


Vérité
terrain

OCR

recouvrement

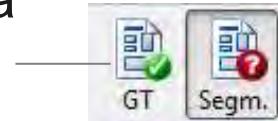
La différentiation des erreurs
s'appuie sur l'ordre de lecture



KO

Évaluation de la segmentation – En pratique

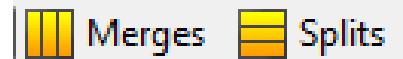
1. Activer l'affichage de la segmentation OCR



2. Chercher les erreurs (par code couleur dans le mode "Overview")



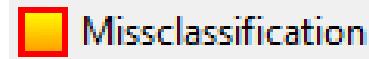
3. Visualiser les erreurs par catégorie (en surbrillance orange)



Merges



Splits



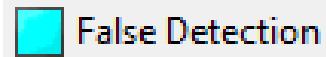
Missclassification



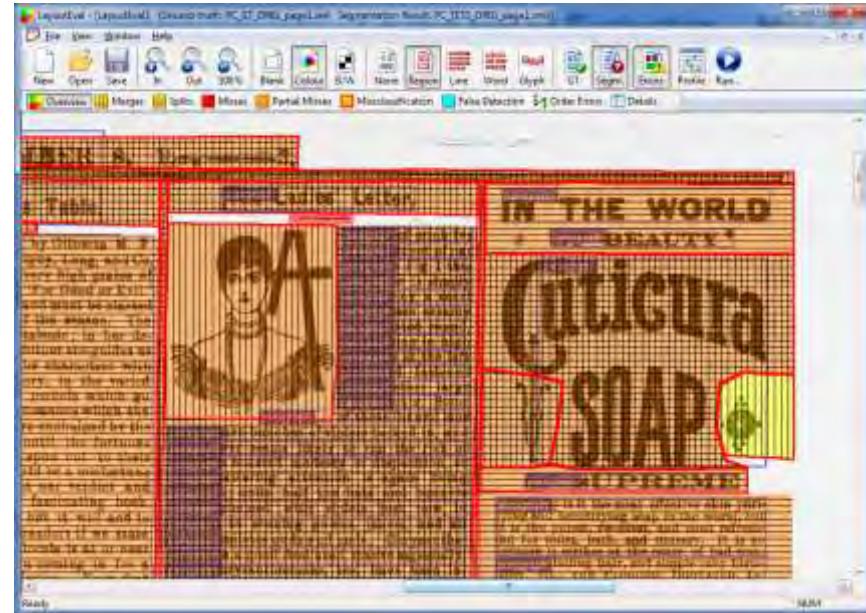
Misses



Partial Misses

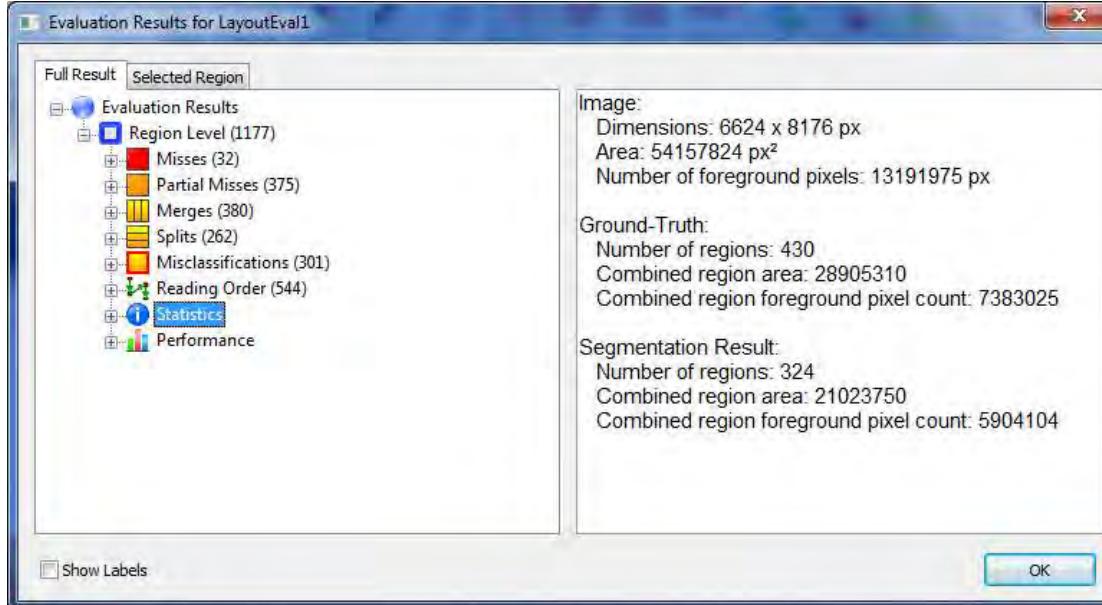


False Detection



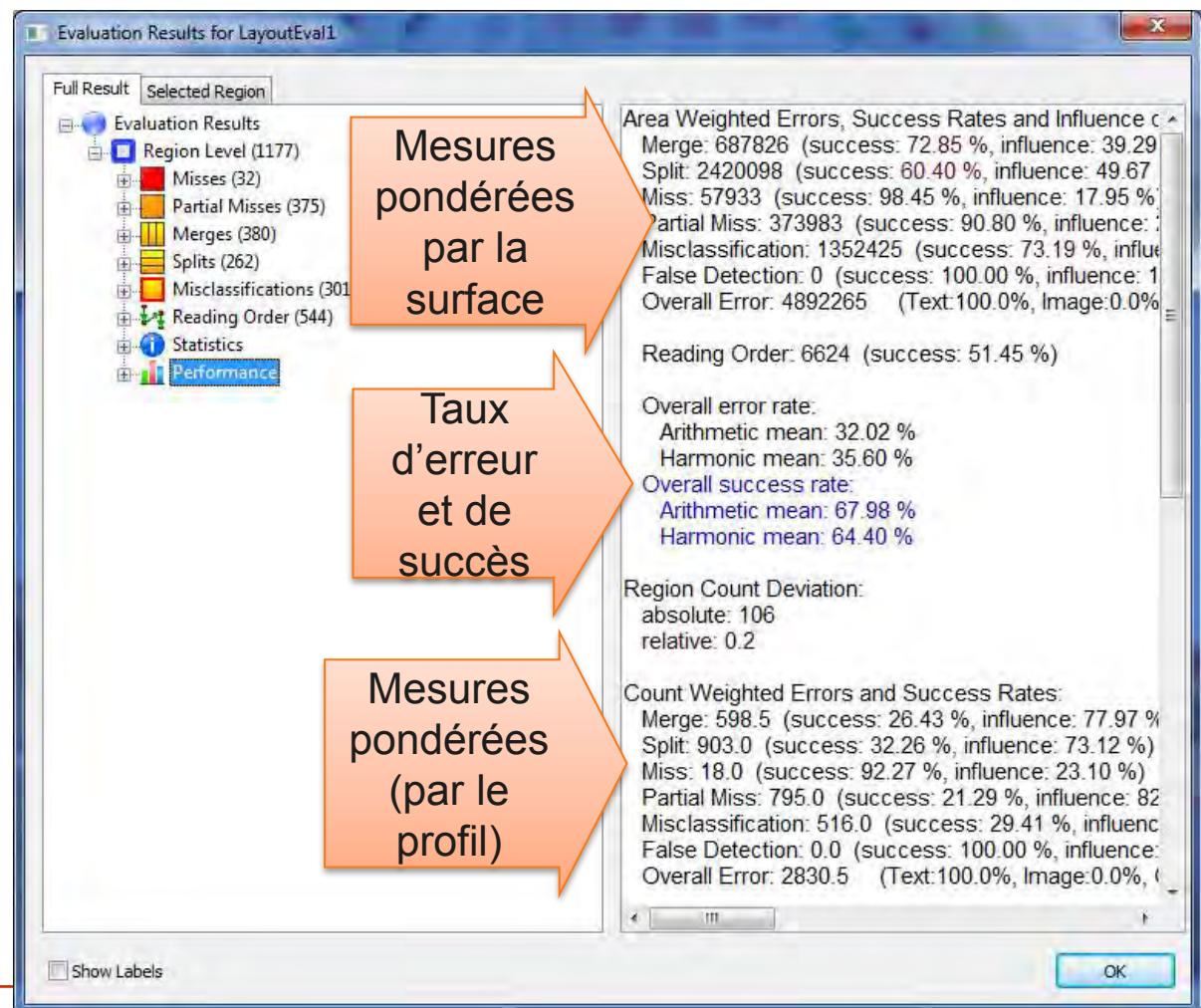
Évaluation de la segmentation – En pratique

1. Clic sur le bouton “Details” —————  Details
2. Clic sur “Statistics”



Évaluation de la segmentation – En pratique

1. Clic sur "Performance"
2. Consulter les taux d'erreur et de succès (en bleu les plus importants)



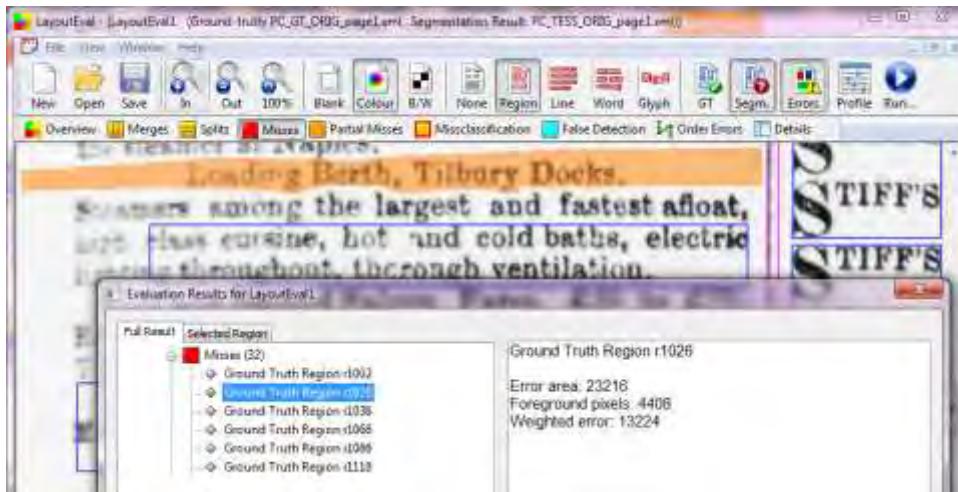
Évaluation de la segmentation – En pratique

1. Activer la vue Vérité terrain (GT)



2. Analyser les résultats en détail, par ex. les omissions :

1. Clic sur “Misses”
2. Clic sur un élément de résultat.

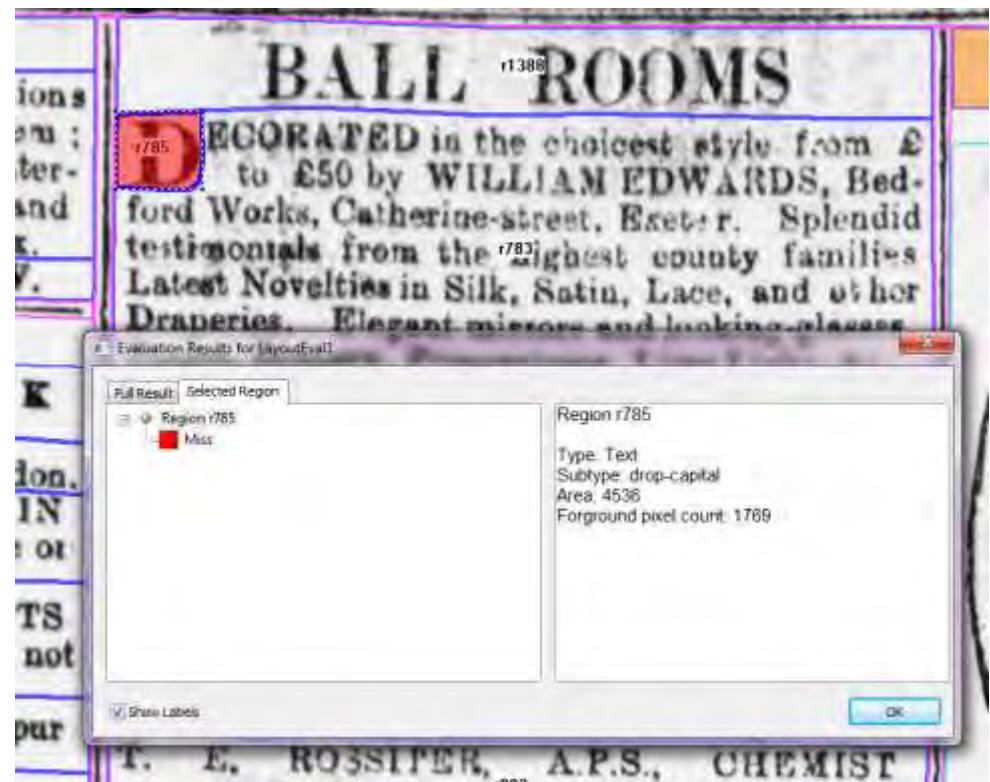
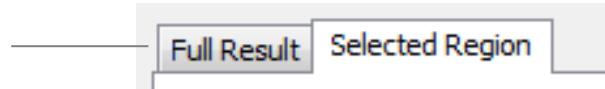


Quantification

Ground Truth Region r1026
Error area: 23216
Foreground pixels: 4408
Weighted error: 13224

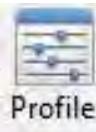
Évaluation de la segmentation – En pratique

1. Sélectionner l'onglet “Selected region”
2. Sélectionner différentes regions : les types d'erreurs sont affichés



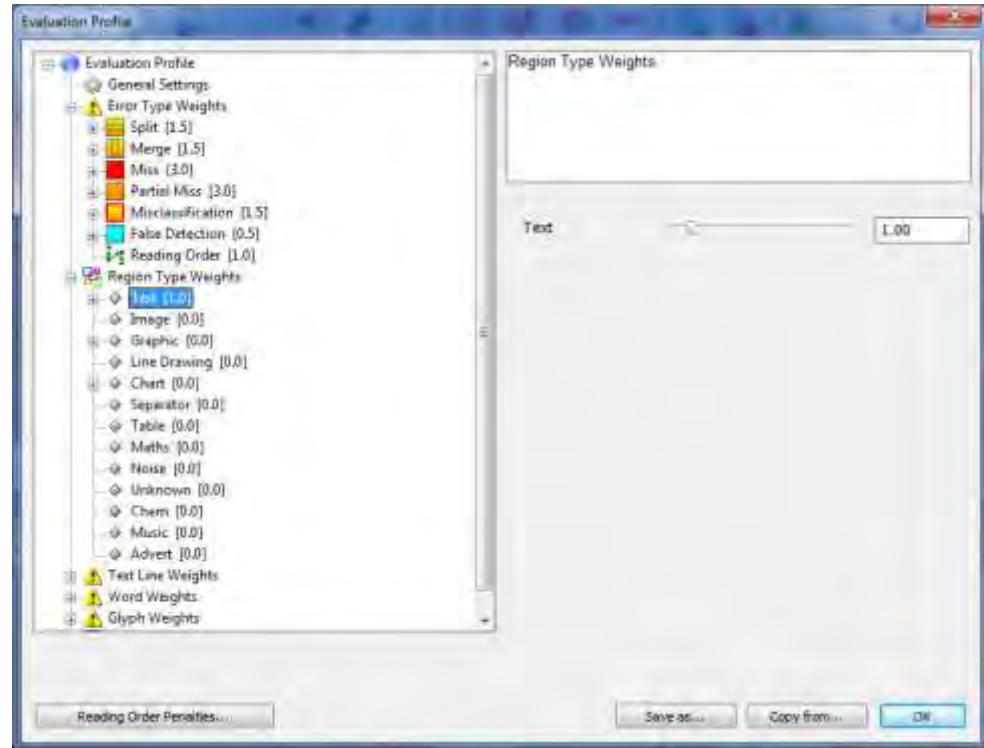
Évaluation de la segmentation – En pratique

1. Clic sur “Profile”
2. “Region Type Weights” donne les poids associés aux natures de contenus
 - Les poids des zones non textuels sont à 0 → les erreurs sur ces régions seront ignorées
 - Même principe pour les catégories d'erreur (“Error Type Weights”)



Profils d'évaluation

- Paramétriser une évaluation pour un cas d'usage donné





europeana
newspapers

Partie 6 : Analyse des performances

Méthodes, outils

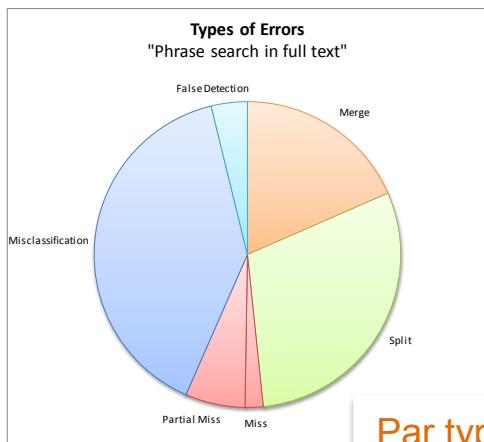
Analyse des performances – Automatisation

- Pour une campagne d'évaluation, nécessité de traiter de nombreux fichiers.
- Les outils sont disponibles en version ligne de commande.
- Traitement par boucle et sortie des résultats au format (CSV).

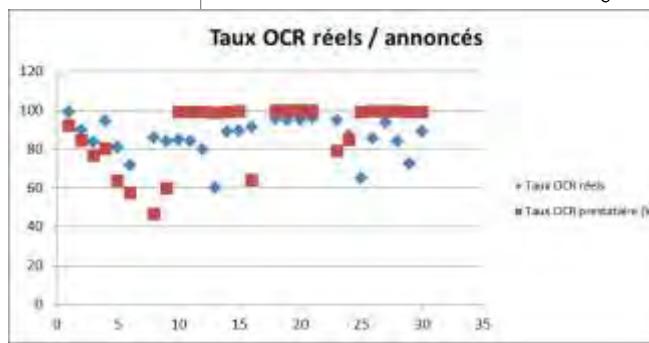
```
C:\Windows\system32\cmd.exe
{folder}.bat"
1
2
C:\Users\clc\Desktop\EMP_Final_Technical_Workshop\Tools\Text Eval 1-2>"evaluate
{folder}.bat"
groundTruth,result,wordIndexMissErrorRate,wordIndexFalseDetectionErrorRate,wordI
ndexSuccessRate,numberOfUniqueWordsInGroundTruth,numberOfUniqueWordsInResult,wor
dCountMissErrorRate,wordCountFalseDetectionErrorRate,wordCountSuccessRate,wordCo
untPrecision,wordCountRecall,wordCountFMeasure,numberOfWordsInGroundTruth,numberOf
WordsInResult,charsInGroundTruth,charsInResult,characterAccuracy,wordsInGround
TruthExclStopWords,wordsInResultExclStopWords,wordAccuracyExclStopWords,wordsInG
roundTruth,wordsInResult,wordAccuracy
input\gt\01.txt,input\res\01.txt,0.222222222222227,0.333333333333333,0.6956521
739130436,11,12,0.318181818181823,0.42857142857142855,0.621761658031089,0.5714
285714285714,0.6153846153846154,0.5925925925925927,13,14,68,77,0.691176470588235
3,13,14,0,46153846153846156,13,14,0,46153846153846156
1
input\gt\02.txt,input\res\02.txt,0.4285714285714286,0.42857142857142855,0.571429
5714285714,2,7,0.4285714285714286,0.42857142857142855,0.5714285714285714,0.57142
85714285714,0.5714285714285714,0.5714285714285714,7,7,57,56,0.9298245614035088,7
,7,0.5714285714285714,2,7,0.5714285714285714
2
C:\Users\clc\Desktop\EMP_Final_Technical_Workshop\Tools\Text Eval 1-2>
```

Analyse des performances

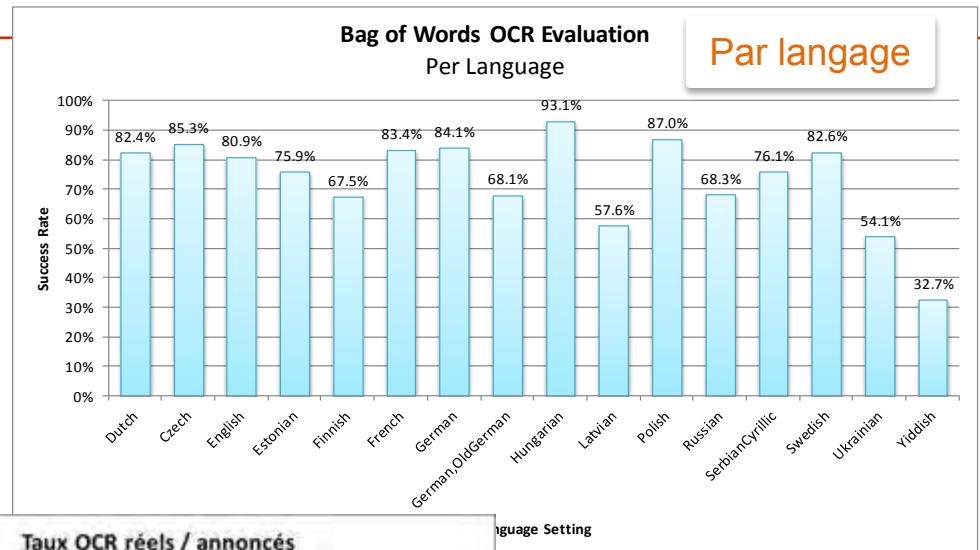
- Import dans un tableau des résultats CSV
- Création de moyennes, graphes, etc.



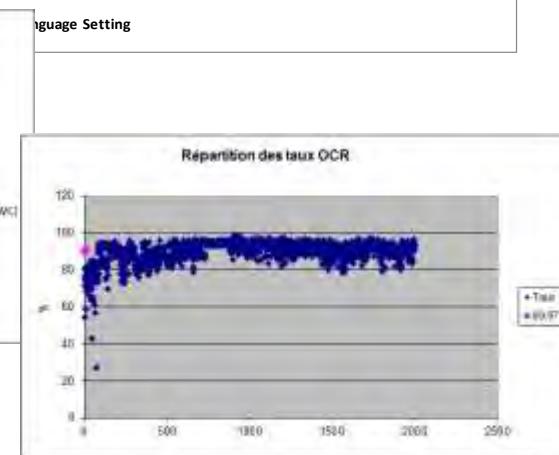
Par types d'erreur



Par taux OCR



Par langage



Prise de décision – Étude avant-projet

Selon les résultats de l'analyse de performance,
il peut être décidé :

- De **lancer** le projet de numérisation.
- De **sélectionner un sous-ensemble** du corpus pour favoriser la qualité OCR.
- D'**améliorer** le processus de numérisation.
- De choisir une **combinaison** des deux options précédentes.
- D'**abandonner** la numérisation du corpus.

Prise de décision – Contrôle qualité en production

Selon les résultats de l'analyse de qualité,
il peut être décidé :

- De demander une **amélioration** de qualité au prestataire ciblée sur les failles identifiées.
- D'**auditer** le processus de numérisation du prestataire afin de l'accompagner dans une démarche d'amélioration de la qualité.
- De **rejeter** les documents non conformes à la qualité attendue de la prestation.
- De **dénoncer** la prestation.



Contact

Christian Clausner (USAL) / c.clausner@primaresearch.org

Jean-Philippe Moreux (BnF) / jean-philippe.moreux@bnf.fr