

Visual perception of unitary elements for layout analysis of unconstrained documents in heterogeneous databases

ICFHR 2014

September 2, 2014

Baptiste Poirriez

Aurélie Lemaitre

Bertrand Couasnon

Irisa / Insa – Intuidoc research group

couasnon@irisa.fr

- **Document Layout Analysis**
 - ◆ **Unconstrained and Heterogeneous Documents**
 - ◆ **Heterogeneous Databases**
 - ◆ **No a priori on the document**

- **Maurdor Campaign dataset**
 - ◆ **2 International Competition in 2013 (www.maurdor-campaign.org)**
 - ◆ **Founded by the French Ministry of Defense**
 - ◆ **Publicly available in 2014**

- **10 000 scanned documents, fully annotated**
 - ◆ **Heterogeneous documents**
 - ✦ **Forms, Tables, Business documents, Correspondences, Faxes
Diagrams, Drawings...**
 - ◆ **French (50%), English (25%), Arabic (25%), Mixed**
 - ◆ **Printed, Handwritten, Mixed**
 - ✦ **1000 writers**

- **Unconstrained and Heterogeneous Documents**
 - ◆ **Methods usually need some homogeneity inside of the collection of documents**
 - ✦ Stability of some blocks of text for example
- **Heterogeneous Databases**
 - ◆ **Methods usually rely on document classification**
 - ◆ **New class for each new type of document**
- **Mixing both is complex**
 - ◆ **No a priori on the document**
 - ◆ **No a priori on the database**

 - ◆ **Only very general a priori on the document elements**

■ Analysis of Salient Elements

◆ Use of Perceptive Vision Mechanisms

- ✦ Already used on homogeneous databases

◆ Perceptive Vision Principle

- ✦ Some contents are salient for the human vision in a document
 - Often strongly structuring
- ✦ Combining several points of view
 - Prediction/verification mechanism
 - Layout is predicted in a global vision of a document and verified with details.

◆ Multi level description of salient elements

■ Primitives at different resolutions

◆ Line segments

◆ Connected components

◆ Words recognized by OCR (Abbyy FineReader)

- **Detection of orientation**
 - ◆ Perceptive vision
 - ◆ Detection at low resolution of text lines (line segments)
 - ◆ Detection of the main direction of writings

- **Iteratively find the most structuring and salient elements**
 - ◆ 1: Tables, separators, boxes
 - ◆ 2: Latin printed blocks of text
 - ◆ **3: Remove of the detected salient elements : New Segmentation**
 - ◆ 4: Handwritten and Arabic text lines
 - ◆ 5: Graphic regions

- **Description of each salient element**
 - ◆ General knowledge on each element

■ DMOS-P: a generic perceptive method

◆ EPF language

- ◆ Grammatical description of a document
- ◆ Generation of an adapted recognition system

◆ Validated

- ◆ Wide range of kind of documents
- ◆ More than 700,000 pages processed
- ◆ On homogeneous documents

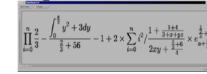
■ Application of DMOS-P for Heterogeneous Documents

- ◆ Description with EPF of salient elements
- ◆ New Segmentation during parsing
- ◆ Combine different levels of perception
 - ◆ Perceptive Layers
- ◆ Combine different salient elements

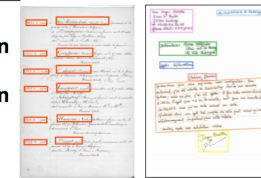
■ Optical Music Recognition



■ Mathematical Formulae Recognition



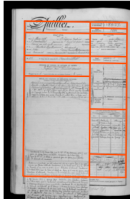
■ Handwritten Structure Recognition



■ Tennis Court Recognition



■ Old Forms Recognition



■ Hierarchical Table Structure Recognition



■ Handwritten Flow Chart Recognition



1: Tables, separators, boxes

■ Perceptive vision

- ◆ Line segments at different resolution
- ◆ Rulings detection

■ General description of a table

- ◆ 2 crossing rulings
- ◆ Table frame on the cross
 - ◆ Virtual rulings
- ◆ Set of rows
- ◆ Set of columns
- ◆ Recursive table detection inside each cell

■ General description of

- ◆ Boxes
- ◆ Separators

Quantité	Article	Unités	Description	Remise %	I.H.T.	Prix unitaire	Total
4	R19	20	Stylos	—		4 €	80€
17	T73	5	Colle VHO	—		8 €	96€
74	866	10	Gammes <i>blanc</i>	—		2,90€	214,60
100	U72	1	Compas	—		6 €	600
101	C88	1	Equeres	—		6,5€	6565
Sous-total							1047,10
I.P.S.							
I.V.O.							
Livraison							
Montant à verser							1047,10

Quantité	Article	Unités	Description	Remise %	I.H.T.	Prix unitaire	Total
4	R19	20	Stylos	—		4 €	80€
17	T73	5	Colle VHO	—		8 €	96€
74	866	10	Gammes <i>blanc</i>	—		2,90€	214,60
100	U72	1	Compas	—		6 €	600
101	C88	1	Equeres	—		6,5€	6565
Sous-total							1047,10
I.P.S.							
I.V.O.							
Livraison							
Montant à verser							1047,10

2: Latin printed block of text

- **Words detected by OCR (Abbyy Finereader)**
 - ◆ **On the complete document**
 - ✦ Errors on mixed (printed/handwritten) documents
 - ◆ **Added in a perceptive layer**

- **General description of Latin Printed block of text**
 - ◆ **Blocks built on words with a high level of confidence**
 - ◆ **Text lines and blocks**
 - ✦ With constraints on size, alignments (left, right, centered), text columns
 - ✦ Very general knowledge on blocks of text

3: New Segmentation (Remove of Table Rulings, Separators, Latin Printed Text)

cerfa N° 20-3265

REPUBLIQUE FRANCAISE Modèle n°10

DÉCLARATION D'ACQUISITION, VENTE, CESSION OU MISE EN POSSESSION DES ARMES DE 5^e CATÉGORIE II OU DE 7^e CATÉGORIE (1)

(Application des articles 47 et 69 du décret du 6 mai 1995)

IMPORTANT : Quiconque se sera fait délivrer indûment ou aura tenté de se faire délivrer indûment un document administratif, soit en faisant de fausses déclarations, soit en prenant un faux nom ou une fausse qualité, soit en fournissant de faux renseignements, certificats ou attestations, sera puni d'un emprisonnement et d'une amende (cf. article 441-6 et 441-7 du nouveau code pénal). Le demandeur est informé que les renseignements qu'il doit fournir pour satisfaire sa demande, sont mémorisés dans un mode de traitement automatisé. Ces informations seront accessibles aux services de l'Etat compétents pour la réglementation des armes et des munitions et aux services de police et de gendarmerie dans le cadre de leurs attributions légales. Le droit d'accès et de rectification aux informations s'exercera auprès de la préfecture (articles 27 et 34 de la loi du 6 janvier 1978 - article 6 de l'arrêté du 12 mars 1986).

Je soussigné : DUPONT Théo

ACQUÉREUR OU PERSONNE MISE EN POSSESSION	VENDEUR OU CÉDANT
Propriétaire ou détenteur	Propriétaire ou détenteur
Nom : <u>DUPONT</u>	Nom : <u>COLT</u>
Prénoms : <u>Théo</u>	Prénoms : <u>Dwayne</u>
Date de naissance : <u>11/11/1995</u>	Date de naissance : <u>10/11/1975</u>
Lieu de naissance : <u>Strasbourg</u>	Lieu de naissance : <u>Chicago</u>
Demeurant à : <u>Molsheim</u>	Demeurant à : <u>Paris</u>
Rue : <u>route du Rhin</u> N° : <u>19</u>	Rue : <u>ave de Rivoli</u> N° : <u>64</u>

Je déclare acquérir, entrer en possession, céder ou vendre(2) l'arme dont les caractéristiques figurent ci-dessous. Je demande la délivrance du récépissé correspondant.

Je certifie sur l'honneur l'exactitude des déclarations portées ci-dessous. A Molsheim le 20/01/12

Signature : Théo Dupont

(1) A établir en deux exemplaires.
(2) Rayer les mentions inutiles.

CARACTÉRISTIQUES DE L'ARME DÉCLARÉE

<p>I - Arme de poing</p> <p>Type (1) : _____ Marque : _____</p> <p>Modèle : _____ N° matricule : _____</p> <p>Calibre : _____ <input type="checkbox"/> Percussion centrale <input type="checkbox"/> Percussion annulaire <input type="checkbox"/> Canon lisse <input type="checkbox"/> Canon rayé</p> <p>Longueur d'arme : <input type="checkbox"/> ≤ 28 cm <input type="checkbox"/> > 28 cm</p> <p><input type="checkbox"/> Arme automatique <input type="checkbox"/> Semi-automatique <input type="checkbox"/> À répétition <input type="checkbox"/> À un coup</p> <p>Catégorie : _____ Paragraphe : _____</p> <p>II - Arme d'épaule</p> <p>Type (2) : <u>carabine</u> Marque : <u>UMAREX</u></p> <p>Modèle : <u>WIGER</u> N° matricule : <u>0426328</u></p> <p>Calibre : <u>4,5</u> <input type="checkbox"/> Percussion centrale <input checked="" type="checkbox"/> Percussion annulaire</p> <p>Nombre de canons : <u>3</u> <input checked="" type="checkbox"/> Canon lisse <input type="checkbox"/> Canon rayé</p>	<p>Longueur canon : <input type="checkbox"/> ≤ 45 cm <input checked="" type="checkbox"/> > 45 cm et ≤ 60 cm <input type="checkbox"/> > 60 cm</p> <p>Longueur de l'arme : <input type="checkbox"/> ≤ 80 cm <input checked="" type="checkbox"/> > 80 cm</p> <p>Système d'alimentation : <input checked="" type="checkbox"/> Automatique <input type="checkbox"/> Semi-automatique</p> <p><input type="checkbox"/> ≥ 3 coups (y compris la chambre) <input type="checkbox"/> ≤ 3 coups (y compris la chambre) <input type="checkbox"/> Magasin ou chargeur amovible</p> <p><input type="checkbox"/> ≥ 5 coups (rechargement à pompe) <input type="checkbox"/> ≤ 5 coups (rechargement à pompe) <input type="checkbox"/> ≥ 10 coups (chargeur seul) <input type="checkbox"/> ≤ 10 coups (chargeur seul)</p> <p><input type="checkbox"/> A répétition <input type="checkbox"/> Un coup par canon</p> <p>Catégorie : <u>3</u> Paragraphe : <u>7</u></p> <p>III - Arme d'épaule ou de poing semi-automatique ou à répétition</p> <p>- Ayant l'apparence d'une arme automatique de guerre (4^e catégorie I paragraphe 9)</p>
--	--

RÉCÉPISSÉ DE DÉCLARATION D'ACQUISITION, VENTE, CESSION OU MISE EN POSSESSION DES ARMES DE 5^e CATÉGORIE II OU DE 7^e CATÉGORIE I

Pièce présentée : Passeport Carte nationale d'identité Carte résident ordinaire Carte résident privilégiée Carte de séjour ressortissant C.E.E. Etrangers autres documents (les préciser) _____

N° : _____

Délivrée le : _____ Jour _____ Mois _____ Année _____

Par : _____

Date de réception de la déclaration : _____ (Cachet) _____ La préfet,

Récépissé remis le : _____

Transmis au préfet le : _____

cerfa N° _____

n°10

Etat : _____

Je soussigné : DUPONT Théo

ACQUÉREUR OU PERSONNE MISE EN POSSESSION	VENDEUR OU CÉDANT
Propriétaire ou détenteur	Propriétaire ou détenteur
Nom : <u>DUPONT</u>	Nom : <u>COLT</u>
Prénoms : <u>Théo</u>	Prénoms : <u>Dwayne</u>
Date de naissance : <u>11/11/1995</u>	Date de naissance : <u>10/11/1975</u>
Lieu de naissance : <u>Strasbourg</u>	Lieu de naissance : <u>Chicago</u>
Demeurant à : <u>Molsheim</u>	Demeurant à : <u>Paris</u>
Rue : <u>route du Rhin</u> N° : <u>19</u>	Rue : <u>ave de Rivoli</u> N° : <u>64</u>

Je déclare acquérir, entrer en possession, céder ou vendre(2) l'arme dont les caractéristiques figurent ci-dessous. Je demande la délivrance du récépissé correspondant.

Je certifie sur l'honneur l'exactitude des déclarations portées ci-dessous. A Molsheim le 20/01/12

Signature : Théo Dupont

CARACTÉRISTIQUES DE L'ARME DÉCLARÉE

<p>I - Arme de poing</p> <p>Type (1) : _____ Marque : _____</p> <p>Modèle : _____ N° matricule : _____</p> <p>Calibre : _____ <input type="checkbox"/> Percussion centrale <input type="checkbox"/> Percussion annulaire <input type="checkbox"/> Canon lisse <input type="checkbox"/> Canon rayé</p> <p>Longueur d'arme : <input type="checkbox"/> ≤ 28 cm <input type="checkbox"/> > 28 cm</p> <p><input type="checkbox"/> Arme automatique <input type="checkbox"/> Semi-automatique <input type="checkbox"/> À répétition <input type="checkbox"/> À un coup</p> <p>Catégorie : _____ Paragraphe : _____</p>	<p>Longueur canon : <input type="checkbox"/> ≤ 45 cm <input checked="" type="checkbox"/> > 45 cm et ≤ 60 cm <input type="checkbox"/> > 60 cm</p> <p>Longueur de l'arme : <input type="checkbox"/> ≤ 80 cm <input checked="" type="checkbox"/> > 80 cm</p> <p>Système d'alimentation : <input checked="" type="checkbox"/> Automatique <input type="checkbox"/> Semi-automatique</p> <p><input type="checkbox"/> ≥ 3 coups (y compris la chambre) <input type="checkbox"/> ≤ 3 coups (y compris la chambre) <input type="checkbox"/> Magasin ou chargeur amovible</p> <p><input type="checkbox"/> ≥ 5 coups (rechargement à pompe) <input type="checkbox"/> ≤ 5 coups (rechargement à pompe) <input type="checkbox"/> ≥ 10 coups (chargeur seul) <input type="checkbox"/> ≤ 10 coups (chargeur seul)</p> <p><input type="checkbox"/> A répétition <input type="checkbox"/> Un coup par canon</p> <p>Catégorie : <u>3</u> Paragraphe : <u>7</u></p> <p>III - Arme d'épaule ou de poing semi-automatique ou à répétition</p> <p>- Ayant l'apparence d'une arme automatique de guerre (4^e catégorie I paragraphe 9)</p>
--	--

RÉCÉPISSÉ DE DÉCLARATION D'ACQUISITION, VENTE, CESSION OU MISE EN POSSESSION DES ARMES DE 5^e CATÉGORIE II OU DE 7^e CATÉGORIE I

Pièce présentée : Passeport Carte nationale d'identité Carte résident ordinaire Carte résident privilégiée Carte de séjour ressortissant C.E.E. Etrangers autres documents (les préciser) _____

N° : _____

Délivrée le : _____ Jour _____ Mois _____ Année _____

Par : _____

Date de réception de la déclaration : _____ (Cachet) _____ La préfet,

Récépissé remis le : _____

Transmis au préfet le : _____

Je soussigné : DUPONT Th

ACQUÉREUR OU PERSONNE

Propriétaire ou

Nom : DUPONT

Prénoms : Thérèse

Date de naissance : 11.11 11.11

Lieu de naissance : Strasbourg

Demeurant à : Mothain

Rue : route du Rhin

Déclare acquérir, entrer en possession
Je demande la délivrance du récépissé
Je certifie sur l'honneur l'exactitude de

- (1) A établir en deux exemplaires.
- (2) Rayer les mentions inutiles.

I - Arme de

DUPONT Thér

DUPONT

Thérèse

11.11 11.11

Strasbourg

Mothain

route du Rhin

entrer en possession,

I - Arme de

■ New Segmentation

- ◆ No more tables, rulings or Printed Latin text
- ◆ Simplify the description of handwritten and Arabic text lines

■ Perceptive Vision

- ◆ Text lines can be seen at low resolution as line segments
- ◆ Confirm the text line at high resolution with regular aligned connected components

■ Description of text lines detects

- ◆ Handwritten text lines : Latin and Arabic
- ◆ Printed text lines: Arabic and some remaining Latin (OCR Errors)

■ Blocks of text made of text lines



Hi Li!

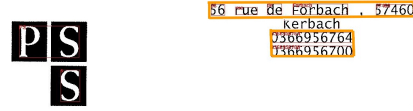
Hey, j'ai appris par Leah que tu venais d'obtenir le job de vendeuse à temps plein par lequel je t'avais fait postuler indirectement, car je savais que tu aimais particulièrement ce magasin. A croire que c'était mon destin de passer devant le magasin du manga "NANA" et qu'il cherchait une vendeuse. Après être venu te voir à Shibuya je suis allée finir mon programme d'étude en chinois à Taipei[...] enfin fini de parler de moi !! J'espère que tu pourras réaliser tous les projets future Dans celui de venir me voir à Paris! La Capitale.

J'espère que tu me "Giras" avant 3 jours et que tu commences ton travail, j'ai hâte de venir te voir dans ce magasin vintage cool. Et n'oublie pas de trouver ton Tatsuya cool.



KUMBAWA
"Socé"

0366956764-0366956700
soncalvesjoc@nuty.com



Hi Li!

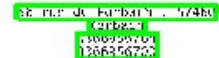
Hey, j'ai appris par Leah que tu venais d'obtenir le job de vendeuse à temps plein par lequel je t'avais fait postuler indirectement, car je savais que tu aimais particulièrement ce magasin. A croire que c'était mon destin de passer devant le magasin du manga "NANA" et qu'il cherchait une vendeuse. Après être venu te voir à Shibuya je suis allée finir mon programme d'étude en chinois à Taipei[...] enfin fini de parler de moi !! J'espère que tu pourras réaliser tous les projets future Dans celui de venir me voir à Paris! La Capitale.

J'espère que tu me "Giras" avant 3 jours et que tu commences ton travail, j'ai hâte de venir te voir dans ce magasin vintage cool. Et n'oublie pas de trouver ton Tatsuya cool.



KUMBAWA
"Socé"

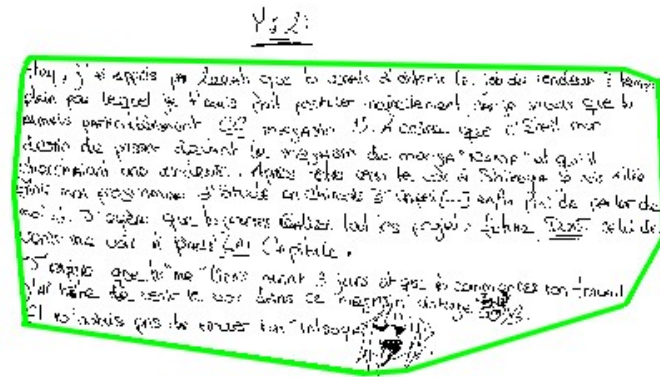
0366956764-0366956700
soncalvesjoc@nuty.com



◆ Words from OCR

◆ Latin Printed Blocks

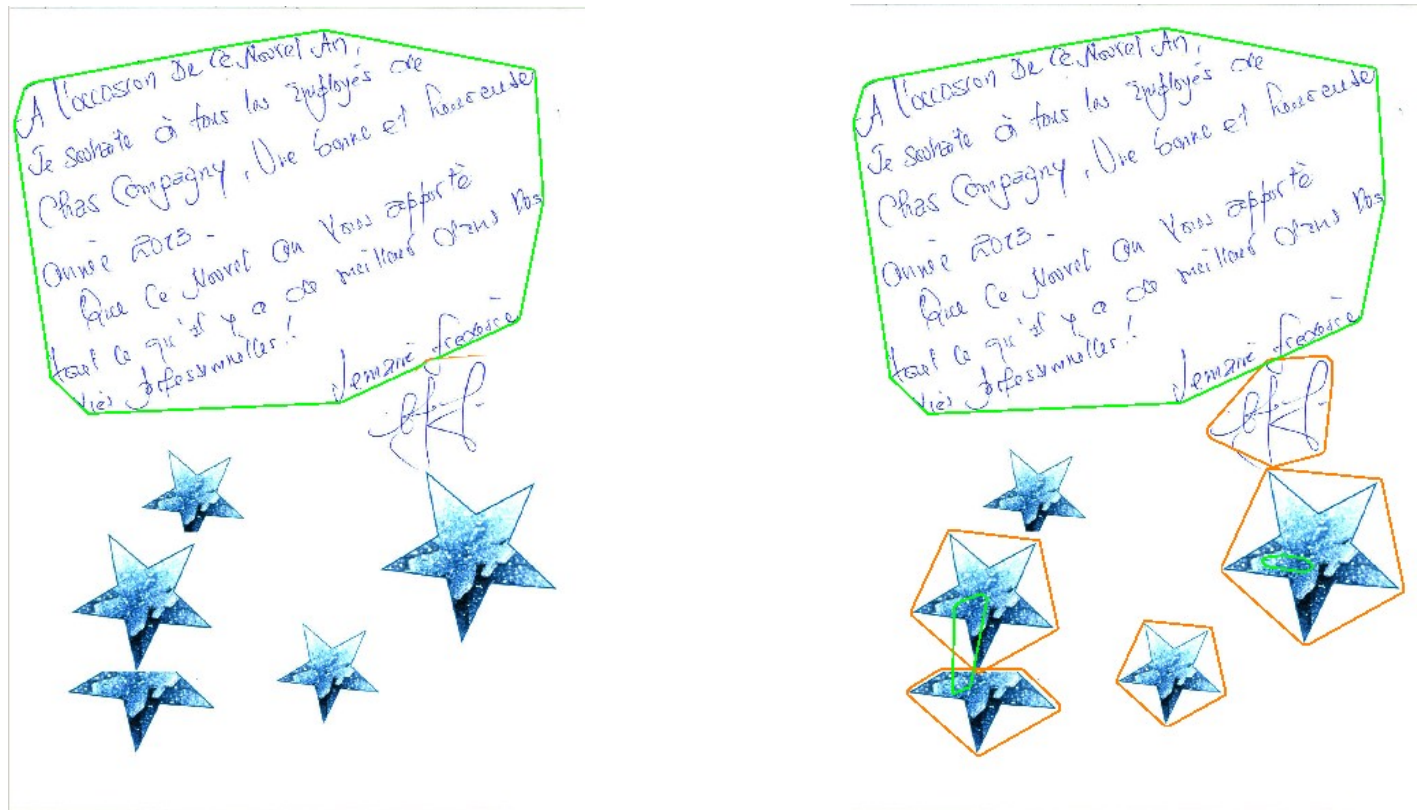
◆ Blocks after Handwritten text lines detection



KUMBAWA
"Socé"

0366956764-0366956700
soncalvesjoc@nuty.com

- Simple description of Graphics
 - ◆ Low resolution
 - ◆ Remaining connected components big enough
- Example



■ Maudor Campaign, Module 1 (Layout Analysis task)

- ◆ Training set : 6,000 documents
- ◆ Dev set : 1,000 documents
- ◆ Test set : 1,000 documents

■ Metrics (see www.maudor-campaign.org for details)

- ◆ ZoneMap (deals with split and merge) : lower better
- ◆ Jaccard (pixel level labeling) : higher better

■ Results

- ◆ 2nd on ZoneMap
- ◆ Close to best on Jaccard

Participant	ZoneMap(%)	Jaccard
Participant 1	48.7	0.45
<i>Our method</i>	59.2	0.44
Participant 2	73.5	0.28

■ Jaccard by class

Participant	Text zone	Graphic zone	Table
Participant 1	0.552	0.394	0.363
<i>Our method</i>	0.553	0.402	0.307
Participant 2	0.307	0.176	0.174

- **DLA on Heterogeneous and Unconstrained Documents**
 - ◆ **Difficult task, open problem**
 - ◆ **General grammatical description of document elements**
 - ✦ Perceptive vision mechanism
 - ◆ **Iterative recognition of salient elements**
 - ◆ **Ability to re-segment the document during the analysis**
 - ◆ **New application of DMOS-P**
 - ✦ Up to now applied on homogeneous databases
- **Improvements**
 - ◆ **Too much split and merge**
 - ✦ Need more homogeneous blocks
 - ✦ Avoid confusion between graphics and text
 - ◆ **Introduce classifiers for adding local information**
 - ✦ Printed/handwritten detection
 - ✦ Language detection
 - ✦ Kind of graphics detection
 - ◆ **Mixing those classifiers with grammatical description**