



# HANDWRITTEN/PRINTED TEXT SEPARATION USING PSEUDO-LINES FOR CONTEXTUAL RE-LABELING

By:

[Ahmad Montaser Awal](#)

Abdel Belaïd

Vincent Poulain d'Andecy

# CONTEXT

- Administrative documents are
  - Noisy
  - Annotated...
- Separation of scripts in administrative documents
  - Annotation extraction
  - Sending each script to a specialized system
  - Noise removal



# STATE OF THE ART

- Printed/handwritten text separation systems share the main steps
  - Preprocessing
    - Removing very small/large connected components
  - Document segmentation
    - Segment the document into basic units
  - Classification
    - Assign each unit to a text class
  - Contextual re-labeling
    - Correct classification errors using neighborhood information

# STATE OF THE ART

## DOCUMENT SEGMENTATION

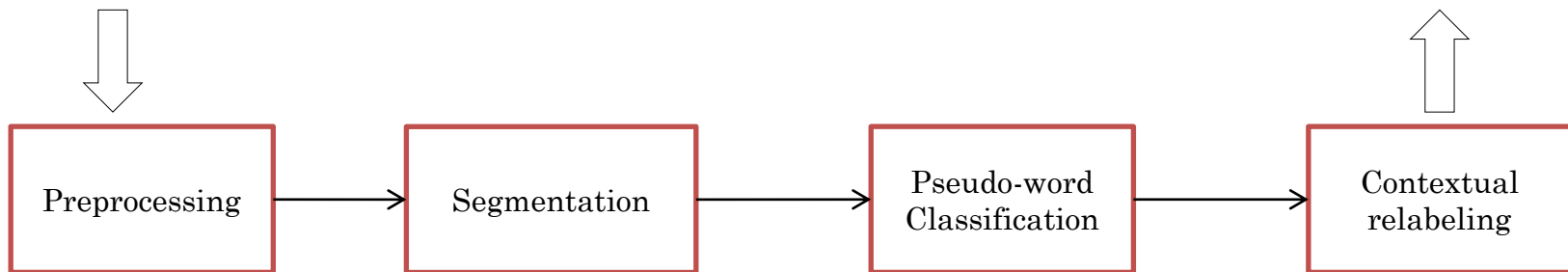
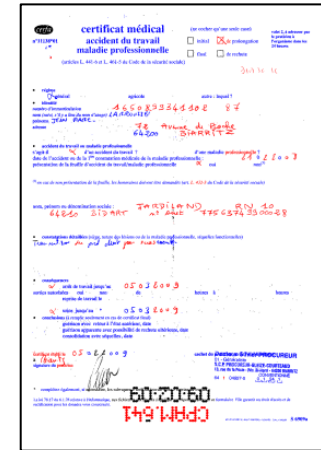
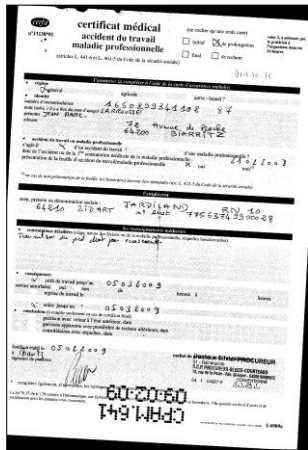
- Text line level (Pal et al. 2001)(Kavallieratou et al. 2004)
  - Lines are assumed to be homogeneous (mono-class)
  - Segmentation using the horizontal projection profiles
- Word level
  - Grouping connected components to approximate words
  - Distance based (Zheng et al. 2004) (Shetty et al. 2007)
  - Morphological operations (Peng et al. 2011) (Zagoris et al. 2014)
- Character level (Fan et al. 1998)
  - Non-cursive scripts (Chinese documents)
  - X-Y cut algorithm

# STATE OF THE ART

## CONTEXTUAL RE-LABELING

- Step1: Define the neighborhood of a given word
  - 4 Nearest Neighbors (Peng et al. 2013) (Zheng et al. 2007)
  - 6 Nearest Neighbors (Shetty et al, 2007)
- Step2: Define criteria to re-label a word based on the labels of its neighborhood
  - Majority voting (kandan et al. 2007)
  - Probabilistic models
    - Markov Random Field (MRF) (Zheng et al. 2007) (Peng et al. 2013)
    - Conditional Random Field (CRF) (Shetty et al. 2007)

# PROPOSED SYSTEM OVERVIEW



\* A. Belaïd, K. Santoch and V. Poulain d'Andecy, "Handwritten and Printed Text Separation in Real Document," *Machine Vision Applications*, vol. 2, 2013

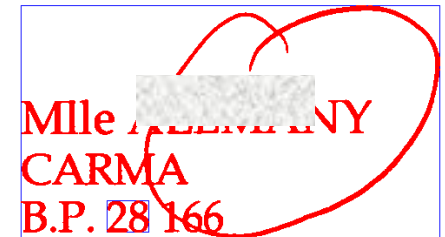
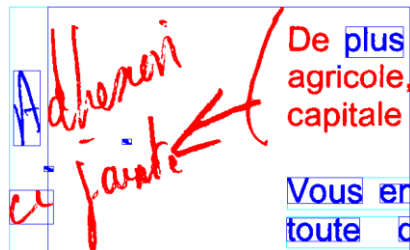
# SEGMENTATION

- Differently from most of existing works, the document is first segmented into pseudo-lines before being segmented into pseudo-words
- Pseudo-line
  - A set of connected components where:
    - Horizontal distances  $< d_H$
    - Vertical distances  $< d_V$
- Pseudo-word
  - A set of connected components belonging to the same pseudo-line
  - Horizontal distance  $< ws$  (word spacing distance estimated automatically for each pseudo-line)



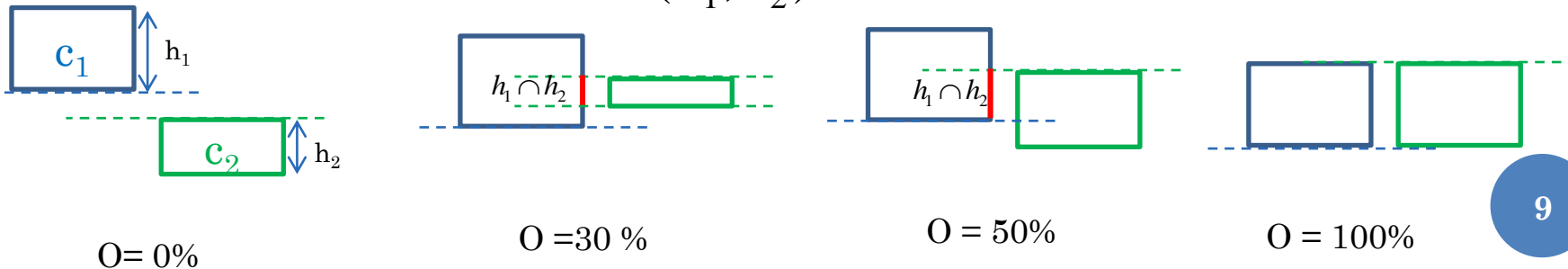
# IMPROVED SEGMENTATION – HEURISTIC

- Avoid vertical connection caused by handwritten annotations

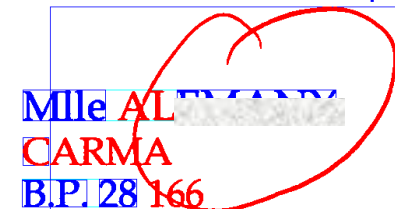
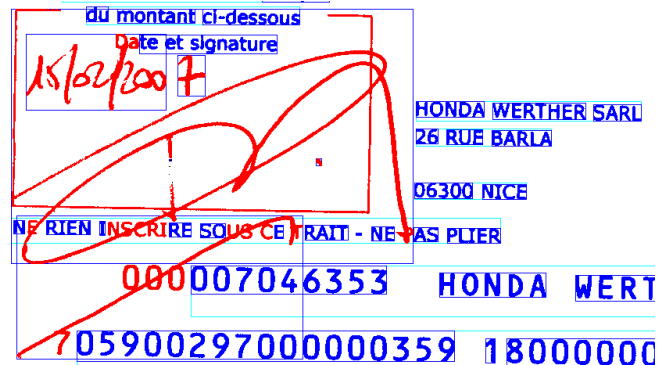
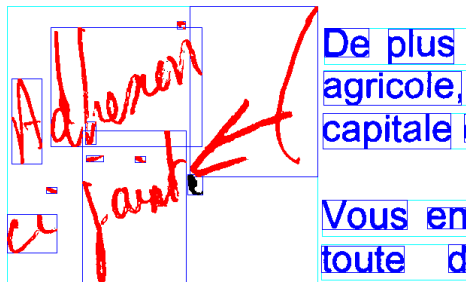
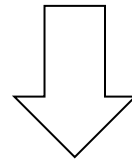
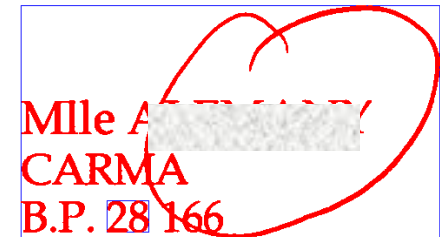
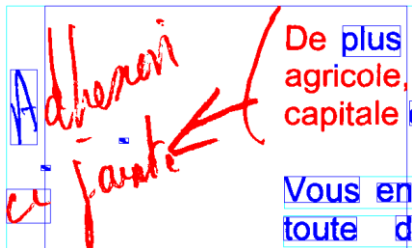


- Use CCs horizontal overlapping

$$o(c_1, c_2) = \frac{h_1 \cap h_2}{\max(h_1, h_2)}$$



# IMPROVED SEGMENTATION – HEURISTIC



# PSEUDO-WORDS CLASSIFICATION

- A pseudo-word is characterized by 137 features
- A multiclass Support vector machines SVM is used to classify a pseudo-word into :
  - Handwritten text
  - Printed text
  - Noise

# CONTEXTUAL RELABELING

- Some classification errors could be corrected using contextual neighborhood
- The label of each pseudo-word is updated based on those of its neighbors
- Local neighborhood
  - K nearest neighbors\*
  - Confidence propagation\*
  - Conditional Random Fields
- Using pseudo-lines
  - Probabilistic model (CRF)
  - Static model

\* A. Belaïd, K. Santoch and V. Poulain d'Andecy, "Handwritten and Printed Text Separation in Real Document," *Machine Vision Applications*, vol. 2, 2013

# CONDITIONAL RANDOM FIELDS (CRF)

- The separation problem can be modeled by CRF
- According to (Nicolas et al. 2007), the probability of a pseudo-word  $w$  is given by:

$$P(X_w | Y_L, Y_C) = \lambda_L f_L + \lambda_C f_C$$

Label field      Local features      Contextual features      Local classifier      Contextual classifier

- Contextual features

- Local classification probabilities of left/right neighbors
- Structural features extracted from the pseudo-word and each neighbor
  - Height ratio
  - Position ratios
  - Density ratio

# RE-LABELING USING PSEUDO-LINES

- Ideally, a pseudo-line represents a text line of the document
- More than 90% of pseudo-lines contain one type of text (printed or handwritten)
- Pseudo-lines define, implicitly, a global horizontal neighborhood relation between the pseudo-words

0.86 0.98 0.80 0.9 0.95 0.99  
Je vous remercie de votre fidélité.

0.97 0.78 0.99 1 0.72 0.91 0.99  
J'attends sur l'honneur

numéro d'immatriculation : 2720259512189 25

# RE-LABELING USING PSEUDO-LINES

- The ***dominant class***  $C_D$  in a pseudo-line is the class with the highest cardinality
- In case of equality of cardinalities, the dominant class is the one with highest average confidence of its pseudo-words
- The label of a pseudo-word is updated:
  - Using a CRF model
  - If it verifies the following condition:

$$(f_i < cf) \vee (|h_i - h_D| < d)$$

Classification Confidence

Certainty factor

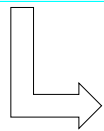
Regularity factor

The diagram shows the mathematical condition  $(f_i < cf) \vee (|h_i - h_D| < d)$  at the top. Three blue arrows point from this condition to three labels below: 'Classification Confidence' (pointing to  $f_i$ ), 'Certainty factor' (pointing to  $cf$ ), and 'Regularity factor' (pointing to  $|h_i - h_D|$ ).

# RE-LABELING USING PSEUDO-LINES

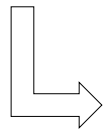
## EXAMPLES

0,74 0,58 0,94 0,75 0,9 0,91 0,96  
 P S : els ont fait faire



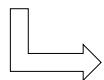
P S : els ont fait faire

0,97 1 0,87 0,88 1 0,94 0,99 0,5 0,5 0,99 1 0,99  
 numéro d'immatriculation : [redacted] 9 5 1 2 1 8 9 2 5



numéro d'immatriculation : [redacted] 9 5 1 2 1 8 9 2 5

0,79 0,73 0,94 0,92 0,99 0,93 0,98 0,92 1  
 Veuillez trouver ci-joint un chèque de = 26,50 € =



No  
Change

0,9 1  
 6 SINGLES → 6 SINGLES

■ Handwritten ■ Printed ■ Noise



# EXPERIMENTATION

## ○ Evaluation

- Pixel level  $\text{pixRate} = \frac{\text{pixels correctly recognised}}{\text{total number of pixels}}$
- Pseudo-word level  $\text{pwRate} = \frac{\text{pseudo- words correctly recognized}}{\text{total number of pseudo- words}}$

## ○ Documents

- Training DB
  - 107 documents (32706 pseudo-words)
    - H: 5888; P: 18078; N: 8740
- Test DB
  - 202 documents (82142 pseudo-words)
    - H: 11970; P: 43705; N: 25190
- All documents are labeled at the pixel level

# RESULTS (1/2)

	System	H%	P%	N%
Previously proposed system*	Proposed system without contextual re-labeling	97.7	96.5	94.3
	k-NN	95.5	97.5	92.3
	Confidence propagation	97.8	96.6	94.0
	CRF	98.5	97.1	94.2
New relabeling methods	Pseudo-lines (CRF): Probabilistic	98.9	97.5	93.5
	Pseudo-lines: Deterministic	98.3	99.2	87.9
Improved segmentation	Pseudo-lines: Deterministic	99.1	99.2	90.1

\* A. Belaïd, K. Santoch and V. Poulain d'Andecy, "Handwritten and Printed Text Separation in Real Document," *Machine Vision Applications*, vol. 2, 2013

# RESULTS (2/2)

System	Docs	<i>pwRate</i>			<i>pixRate</i>			
		H%	P%	ALL%	H%	P%	N%	ALL%
[kandan et al. 2007]	150	-	-	93.2	-	-	-	-
[Zheng et al. 2004]	94	93.0	98,0	98.1	-	-	-	-
[Peng et al. 2013]	82	93.8	95,7	95.5	-	-	-	-
[Shetty et al. 2007]	27	-	-	-	94.8	98.4	89.8	95.7
[Hamrouni et al. 2014]	32	-	-	-	80.0	92.8	-	90.1
Proposed system	202	97.3	99.5	98.7	99.1	99.2	90.1	96.8

# CONCLUSION AND PERSPECTIVES

- Distance based segmentation is not always enough to obtain ‘good’ pseudo-words
  - Heuristics could improve and solve some segmentation problems
- A better performance using pseudo-line based contextual relabeling
- A very good performance compared to the state of the art systems
- In future work:
  - Feature selection
  - Ambiguity layer

Thank *you*