

Fully Convolutional Networks for physical and logical structure analysis of historical newspapers

Yann Soullard, Pierrick Tranouez, Stéphane Nicolas,
Clément Chatelain, Thierry Paquet

University of Rouen Normandy, INSA Rouen Normandy

June 6, 2019



Introduction

State of the art

Fully convolutional networks

Preliminary results

Conclusion

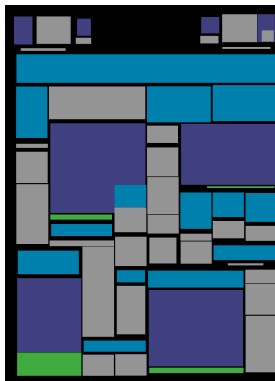
Historical newspapers analysis:

- Field identification in a page
 - ➡ title level 1, 2,..., text, advertisement, image, caption, ...
- Structure analysis
 - ➡ articles, sections, reading order.
- Information extraction
 - ➡ article classification, name entities retrieval,...

Introduction

Field identification in a page, it's difficult!

- Class imbalance
- Not always rectangular blocks
- Various scales
- Various fonts and printed styles
- Confusing classes



Introduction

Field identification in a page, it's difficult!

- Class imbalance
- Not always rectangular blocks
- Various scales
- Various fonts and printed styles
- Confusing classes



Introduction

Field identification in a page, it's difficult!

- Class imbalance
- Not always rectangular blocks
- Various scales
- Various fonts and printed styles
- Confusing classes



State of the art

Newspaper analysis

- HOG features in input of a random forest classifier, rectangular bounding boxes [Lang et al., 2018]
- Contour classification and morphological operations [Vasilopoulos and Kavallieratou, 2017]

Related domains

- Localization of handwritten fields [Lemaitre et al., 2018]
 - ➡ rule-based system
- dhSegment [Oliveira et al., 2018]
 - ➡ Fully Convolutional Network (from a ResNet-50)
- Tools and software (PIVAJ [Tranouez et al., 2015], DMOS [Couasnon and Lemaitre, 2017])

State of the art - PIVAJ

History of PIVAJ

- An offline article analyser and an online viewer and transcription editor
- From a research project to a licensed software suite
- Currently used by National Libraries of Finland and Wales



State of the art - PIVAJ

History of PIVAJ

- An offline article analyser and an online viewer and transcription editor
- From a research project to a licensed software suite
- Currently used by National Libraries of Finland and Wales



State of the art - PIVAJ

PIVAJ offline

1. Pixel level classification (Text, Image, Separators)
2. Connected component, bounding boxes
3. Bottom up structural grid and reading order
4. Text level classification with Random Forests
5. Section and articles assembly through regular expressions

The goal of this current work is twofold:

- Investigate deep learning methods for physical and logical structure analysis,
- Gain in generalization capabilities (being less data-dependent).

Introduction

State of the art

PIVAJ

Fully convolutional networks

FCN

Improvements of convolutional networks

Proposed architectures

Preliminary results

Conclusion

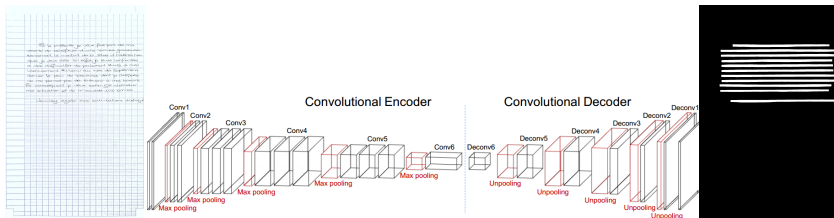
Fully convolutional networks (FCN) were first introduced in semantic segmentation [Long et al. 2014].

FCN are convolutional networks without a dense layer.

Removing dense layers has many advantages:

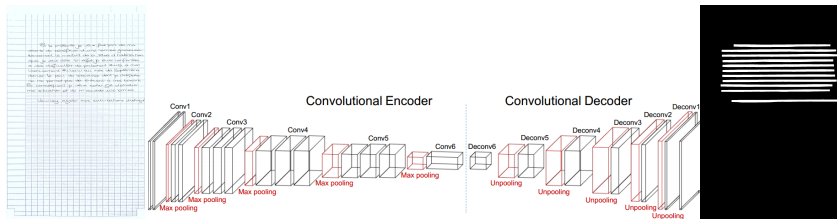
- ▶ the number of parameters is highly reduced,
- ▶ they can deal with variable input size,
- ▶ spatial information is kept inside the network: we can produce a pixel labeling.

Pixel labeling ➡ How to get an output with the same dimension as the input image?

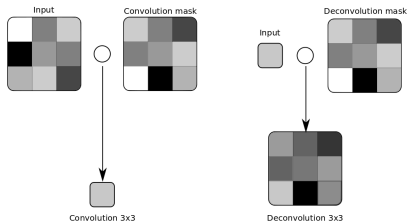


- Transposed convolutions,
- Unpooling layers,
- Dilated convolutions.

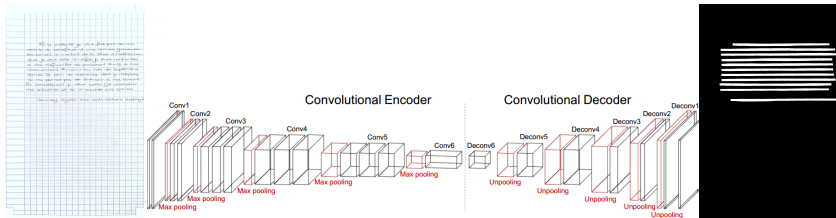
Pixel labeling \Rightarrow How to get an output with the same dimension as the input image?



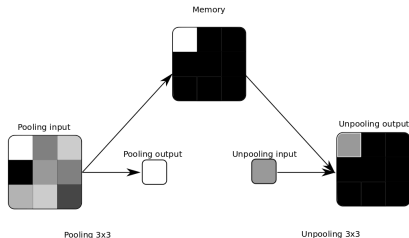
- Transposed convolutions,



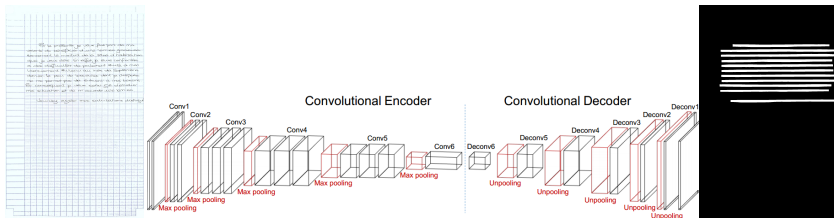
Pixel labeling → How to get an output with the same dimension as the input image?



- Unpooling layers

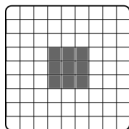


Pixel labeling ➡ How to get an output with the same dimension as the input image?

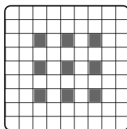


- Dilated convolutions.

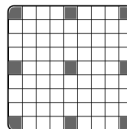
$r = 1$



$r = 2$



$r = 4$

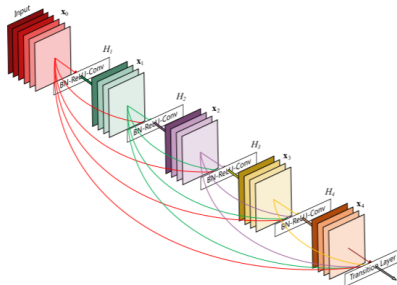


Improvements of convolutional networks

Strengthen the capacities of convolutional networks

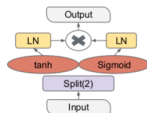
- Dense Nets

➡ Generalize a residual architecture



- Gated Nets

➡ Auto-attention mechanism

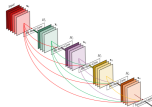


Improvements of convolutional networks

Strengthen the capacities of convolutional networks

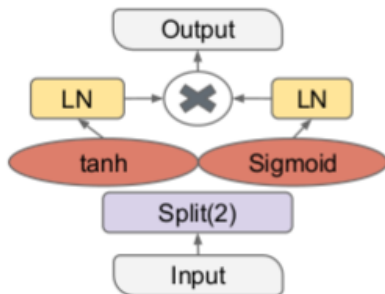
- Dense Nets

- ➡ Generalize a residual architecture



- Gated Nets

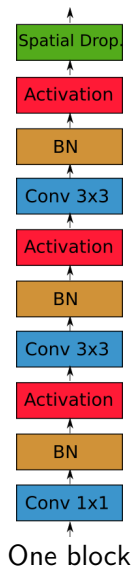
- ➡ Auto-attention mechanism



Proposed architectures

We experiment 3 network architectures:

- **Dilated FCN 22:** FCN with 21 dilated convolutions (7 convolutional blocks with dilation rates 1-2-4-8-4-2-1)
- **Dense Net 22:** the same architecture than the Dilated FCN 22 with a dense architecture between blocs.
- **Gated Net 22:** the same architecture than the Dense Net with a gate at the end of each bloc.



Introduction

State of the art

PIVAJ

Fully convolutional networks

FCN

Improvements of convolutional networks

Proposed architectures

Preliminary results

Conclusion

Data sets

Historical newspapers:

- **Le Matin**

French daily newspaper from 1883 to 1944.
801 examples in training, 9 classes (6 in the experiments)

- **GBNLA**

ICDAR competition on historical newspapers.
66 examples in training from 6 newspapers, 4 classes

- **Luxembourg Historical Newspapers**

980 examples in training from 1877, 11 classes

Input image size are reduced to at least 1200 × 750 by keeping the ratio.



Results

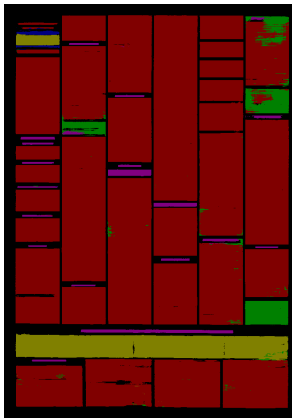
Accuracy (%) on each newspaper, training from scratch.

	Le Matin	GBNLA	Luxembourg
Dilated FCN 22	85.73	87.17	81.55
Dense Net 22	86.31	89.06	81.66
Gated Net 22	85.19	86.57	82.23

- Dense architecture improves performance,
- Having half the feature maps for the gates can be beneficial,
- Gated networks are prone to overfitting.

Additional observations (1/3)

Doubt is generally legitimate



du style. M. Auguste Vaquerie aurait beaucoup aimé ce drame où les revendications de l'esprit et de la chair sont exprimées en alexandrins sonores et bien coupés pour le théâtre; cette application du romantisme du style à la défense de l'indépendance de la femme et du cultes des prêtres a été très applaudie.

Point de réserves à faire sur l'interprétation. L'apostrophe a été très bien jouée par Mme Deubrive, Mlle Vergeuil, de l'Odéon; MM. Prad et Dorsey. — H. G.

On écrit de Bruxelles que M. Van Dyck vient de signer un engagement avec la direction du théâtre de la Monnaie, de Bruxelles, pour aller donner, à la fin de la saison prochaine, sept représentations de *Tannhäuser* et de *Lohengrin*.

Au théâtre du Farc, M. Leberty, qui s'est fait remarquer dans *Cabotine*, a renouvelé son engagement avec M. Alhambra.

A l'Opéra, l'examen trimestriel de la danse, qui devait avoir lieu vers le milieu du mois de juillet, se trouve avancé pour raison de service; il aura lieu le 27 du courant.

La reprise d'Aïda n'aura lieu que dans les premiers jours de juillet. Ce n'est pas M. Alvarez, retenu à Londres, qui chantera Elvira; mais, Mme M. Dupuyroux, qui chantera Elvira, rentrera qu'à partir du 1^{er} août, étant allée à Londres jusqu'à cette date.

C'est irrévocablement le samedi 22 juin qu'aura lieu, au Nouveau-Théâtre, 15, rue Blanche, à huit heures du soir, le huitième spectacle du théâtre de l'Œuvre. Au programme : *Brand*, d'Offenbach.

La veille, vendredi 21, à huit heures du soir, répétition générale.

En raison de la longueur du spectacle, la représentation commencera très exactement aux heures indiquées ci-dessus et, durant les entrées, les portes donnant accès dans la salle seront tenues rigoureusement fermées.

On a commencé à l'Ambro-Comique, les répétitions générales du *Proteus* n° 6, dont voici la distribution définitive.

Michel Servan MM. Desrois
De Merville J. Lecot
Louis Disboure Gémier

REVUE DE PHOTOGRAPHIE

Une découverte toute récente vient d'être faite par le docteur Barillet touchant les procédés de l'hydrolyse formolique au formol. C'est en étudiant la capacité antiseptique considérable de cet agent qu'il a été constaté la faculté que possédait le formol de servir la pellicule en la rendant imputrescible. Cette propriété est celle qui intéresse le photographe. Depuis longtemps, pour réduire la pellicule imputrescible après l'usage des hautes températures, on a continué de plonger la plaque pendant quelques minutes dans un bain d'alun. Avec le formol, on obtient un résultat plus complet et plus rapide.

Voici le procédé qui donne les meilleurs résultats : on passe au bain composé de formol du commerce et d'eau distillée à parties égales; immerger complètement le cliché après lavage soigné au sortir du bain d'hyposulfite; après deux ou trois minutes, retirer le cliché et le placer sur un égouttoir pour le séchage. Quelques praticiens conseillent de laver à l'eau bouillante et de sucer au feu.

Notre procédé permet d'obtenir en dix minutes un cliché absolument sec prêt au tirage des épreuves positives ou des agrandissements.

Toutes les communications touchant la photographie doivent être adressées : *Photographie Le Matin*, 25, rue d'Armenville.

Lyon, même Villefranche et Moulins, ont fait des larges tirages dans nos parages.

Cours pratiques :
Vins blancs ordinaires..... 50 à 55 fr.
— chics..... 60 75
— supérieurs..... 80 85
Vins rouges ordinaires..... 60 65
— chics..... 65 70

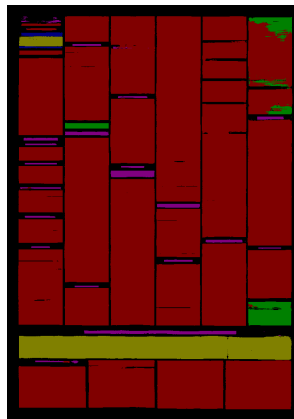
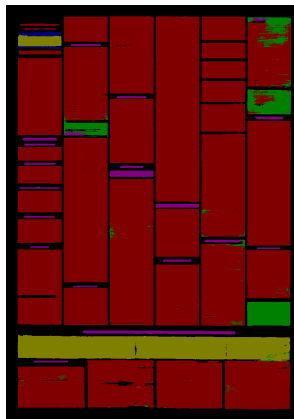
Par le Procureur Dumas et le Secrétaire Estimationnel le docteur Barillet touchant les procédés de l'hydrolyse formolique au formol. C'est en étudiant la capacité antiseptique considérable de cet agent qu'il a été constaté la faculté que possédait le formol de servir la pellicule en la rendant imputrescible. Cette propriété est celle qui intéresse le photographe. Depuis longtemps, pour réduire la pellicule imputrescible après l'usage des hautes températures, on a continué de plonger la plaque pendant quelques minutes dans un bain d'alun. Avec le formol, on obtient un résultat plus complet et plus rapide.

Le Matin newspaper.

- This is observed in case of confusing observations,
- There is sometimes labelling errors and clumsy labelling.

Additional observations (2/3)

Fine-tuning on a specific time period is important

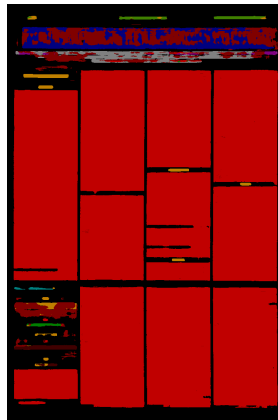
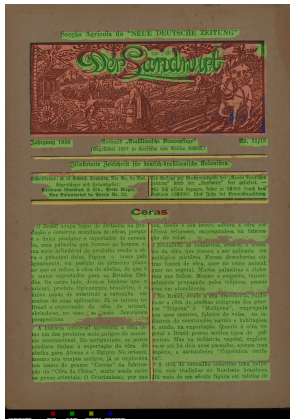


Le Matin newspaper.

- Predicted blocs are more precise and more confident,
- The model tends to overfit (help to label confusing classes).

Additional observations (3/3)

Class imbalance must be taken into account



GBNLA (left) and *Luxembourg* (right) newspapers.

- Text parts are well identified,
- Infrequent classes (e.g. levels of title) are difficult to recognize (even if there is highly similar pages in training).

Introduction

State of the art

PIVAJ

Fully convolutional networks

FCN

Improvements of convolutional networks

Proposed architectures

Preliminary results

Conclusion

Conclusion

We investigate **Fully Convolutional Networks for analyzing historical newspapers**

- FCN can manage input images of variable size,
- FCN provide pixel labeling,
- Recent improvements (e.g. Gate, denseNet) strengthen the model capacities with auto-attention and multiscale analysis.

Regarding **future works**

- Increase the receptive fields \Rightarrow more context in prediction,
- Hierarchical analysis \Rightarrow class imbalance,
- Joint multi-scale analysis (dense-like architecture) with attention mechanism \Rightarrow recognize blocks of various scales,
- Extract polygon boxes \Rightarrow further analysis (removing regions without interest, reading order...).

Thanks for your attention!