

Visual and Textual Feature Fusion for Document Analysis

Patricia Medyna Lauritzen de Lucena Drumond

Supervisor

Prof. Dr. Teófilo Emidio de Campos

Prof. Dr. Fabricio Ataide Braz

Programa de Pós-Graduação em Informática



UnB

Departamento de
Ciência da Computação

11th November 2022

- Introduction
- Background
- Proposal
- Experiments
- Results
- Schedule
- Reference

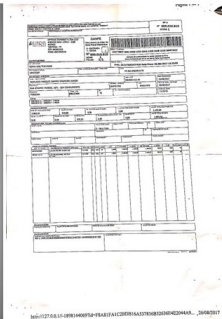
• Motivation



(a)



(b)



(c)



(d)

Document Images

- **Problem**

- Document information extraction
- Wide variety of layout
- Visual and textual features

- **Main Goal**

- Implement and evaluate document processing methods that combine textual information and layout with low computational cost.

• **Specific Goals**

- to propose the joint of textual and layout features for information extraction.
- to evaluate this approach for document classification and page segmentation.
- to compare the models with baselines.

• Contributions

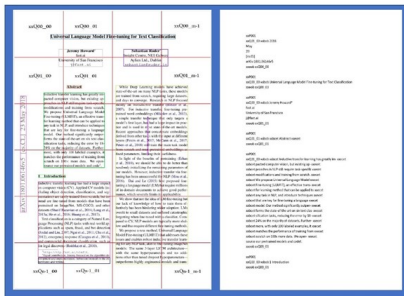
- A novel approach to fuse textual and layout information.
- The simple yet effective model.
- The source code of our library, which is available from <https://github.com/patriciamedyna/LayoutQT>

Table: Document AI: models, modality, backbone and datasets.

Models	Modality	Backbone	Datasets
Asim et al. (2019) [1]	T + I	InceptionV3 Multi-channel CNN	Tobacco-3482 RVL-CDIP
Audebert et al. (2020) [2]	T + L	Multimodal Neural Network	Tobacco-3482 RVL-CDIP
LayoutLM (2020) [8]	T + L	Transformer BERT	FUNSD SROIE RVL-CDIP
Wiedemann and Heyer (2021) [7]	T + I	CNN (VGG16) MLP	Tobacco800 German dataset
Braz et al. (2021) [3]	I only	CNN (VGG16) EfficientNet	Tobacco800 AI.Lab.Splitter
LayoutLMv2 (2021)[9]	T + L + I	Transformer	FUNSD SROIE CORD Kleister-NDA RVL-CDIP DocVQA

LayoutQT

Input



Processing Language Model

Tokenization

Training

Evaluation

Proposal

LayoutQT

Expression	Description
xxbob	begin-of-block
xxeob	end-of-block
xxbcet	begin-of-center-text
xxecet	end-of-center-text
xxQhi_vj	Quadrant tag, line hi and column vj
xxPk	Page tag, number k

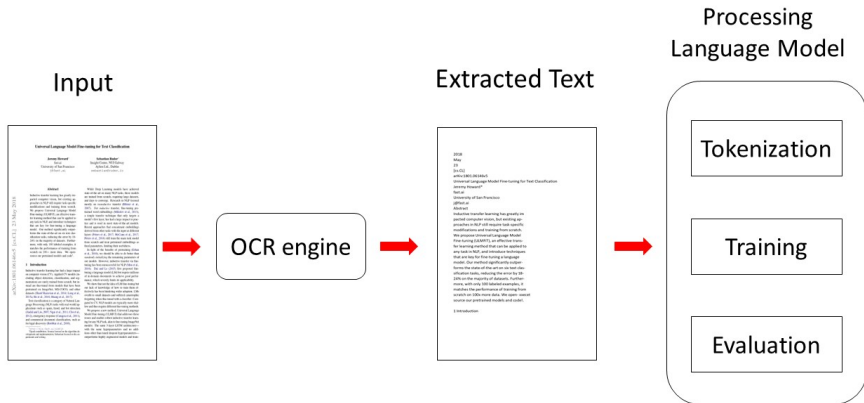
Algorithm 1 LayoutQT Algorithm

Input: multi page document

Output: tokenized text t

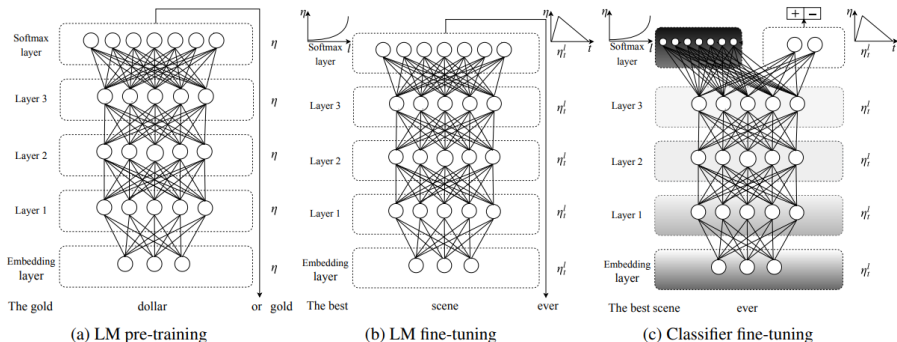
```
1: for  $page = 0, \dots, N - 1$  do
2:    $t+ =$  add page token (where  $+ =$  means insert symbol in string  $t$ )
3:   triage each word bounding boxes into line and group
4:   triage groups into coherent page columns
5:   for each group do
6:      $t+ =$  quadrant coordinate of group top left corner
7:     for each text line in this group do
8:       check line centralization w.r.t. its page column center position
9:       if the line is centralized then
10:         $t+ =$  centre tag
11:       end if
12:        $t+ =$  textual contents of the line
13:       if the line is centralized then
14:         $t+ =$  centre tag
15:       end if
16:     end for
17:      $t+ =$  quadrant coordinate of group bottom right corner
18:   end for
19: end for
```

● Baseline



- **Page Stream Segmentation**
 - LSTM
 - ULMFiT (AWD-LSTM)
 - BERT
- **Document Type Classification**
 - ULMFiT (AWD-LSTM)
 - BERT

Universal Learning Fine-Tuning (ULMFiT)



Source: Howard and Ruder (2018) [5]

- **Dataset:**

- Tobacco-800: 1290 images: 831 training, 200 validation, and 259 test images.
- RVL-CDIP: 400,000 images (16 classes with 25,000 images per class): 320,000 training, 40,000 validation, and 40,000 test images.

● Experimental Settings

- The model is trained with a batch size of 128 and a sequence length of 150 for 100 epochs using NVIDIA Tesla V100 32GB GPU.
- LSTM: 256 nodes fully connected with activation “ReLU” and a dropout of 0.3, binary cross-entropy as a loss function with softmax activation and Adam as an optimizer.
- AWD-LSTM [6]: 3 layers, 1152 hidden sizes and 24M parameters.
- BERT [4]: 12 layers, 768 hidden sizes, 12 self-attention heads and 110M parameters.

Result of binary classification on Tobacco 800 dataset

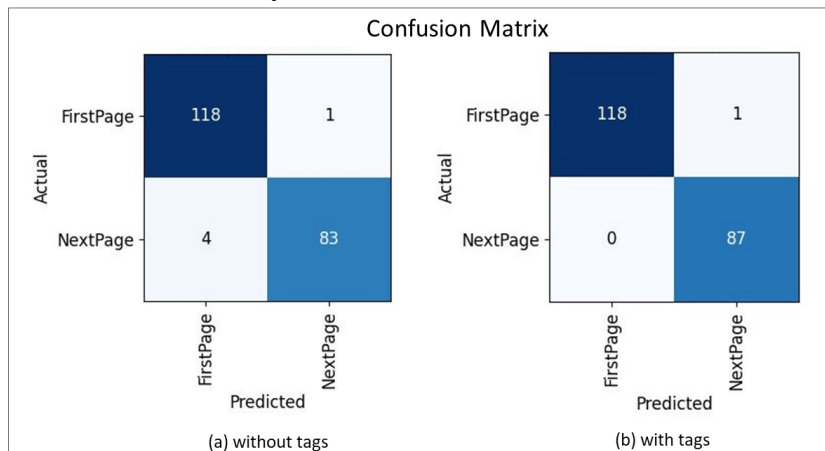


Table: Result of binary classification on Tobacco 800 dataset.

Model	Modality	Backbone	Accuracy	F1-score
Wiedemann et al. (2019) [7]	text + image	VGG16*	91.1%	90.4%
Braz et al. (2021) [3]	only image	VGG16*	92.0%	91.9%
Braz et al. (2021) [3]	only image	EfficientNet-B0*	83.7%	81.9%
Baseline	only text	LSTM	84.1%	82.9%
LayoutQT	text + layout	LSTM	85.9%	86.1%
Baseline	only text	<i>BERT_{BASE}</i>	92.2%	92.0%
LayoutQT	text + layout	<i>BERT_{BASE}</i>	93.0%	93.0%
Baseline	only text	ULMFiT (AWD-LSTM)	97.5%	97.9%
LayoutQT	text + layout	ULMFiT (AWD-LSTM)	99.5%	99.1%

Result of the document types classification on RVL-CDIP dataset

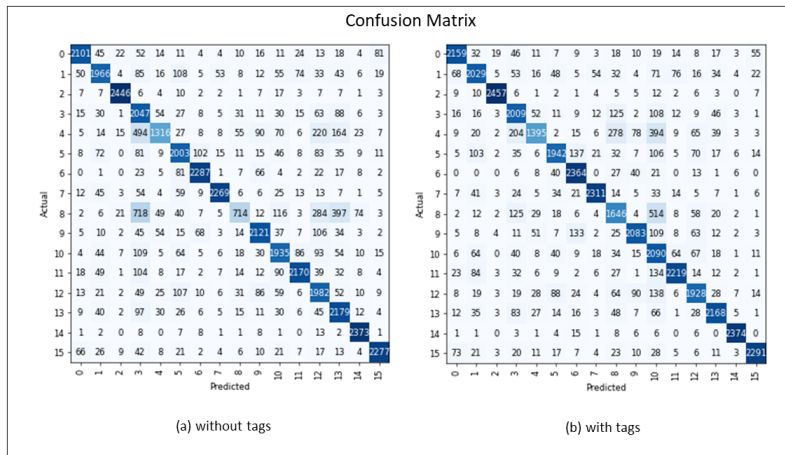


Table: F1-score of the document types classification on RVL-CDIP dataset








Class	Document Type	Baseline AWD-LSTM	LayoutQT AWD-LSTM	Baseline <i>BERT_{BASE}</i>	LayoutQT <i>BERT_{BASE}</i>
0	letter	85.5%	87.8%	83.7%	86.0%
1	form	78.8%	81.3%	77.8%	77.3%
2	email	97.2%	97.6%	93.0%	96.0%
3	handwritten	84.9%	83.3%	63.6%	80.0%
4	advertisement	55.2%	58.5%	66.0%	70.0%
5	scientific report	80.6%	78.2%	74.8%	80.3%
6	scientific publication	89.1%	92.1%	87.4%	89.0%
7	specification	91.9%	93.6%	90.7%	91.0%
8	file folder	31.9%	73.5%	64.0%	73.8%
9	news article	86.4%	84.9%	78.8%	82.6%
10	budget	77.8%	84.0%	78.1%	82.3%
11	invoice	87.9%	89.9%	81.4%	85.9%
12	presentation	79.9%	77.8%	70.3%	81.1%
13	questionnaire	90.0%	89.5%	83.7%	87.9%
14	resume	93.6%	93.7%	98.6%	98.3%
15	memo	91.9%	92.5%	85.4%	90.0%
Average		80.4%	83.6%	80.1%	84.5%

- Engineering Applications of Artificial Intelligence – Qualis A1 (in review) Drumond, P. M. and Leite, L. and Campos, T. and Braz, F. “LayoutQT - Novel Preprocessing Combining Visual and Textual Features”. (2022)

Table: Summary of research activities planning.

Activity	Nov/22	Dec/22	Jan/23	Feb/23
Experiments	✓	✓		
Baseline on VICTOR dataset	✓			
Current model on VICTOR dataset	✓			
Model parameter adjustment		✓	✓	
Training and validation on chosen datasets		✓		
Thesis Writing	✓	✓		
Background Update + Related Works	✓			
Methodology		✓		
Results and Discussion		✓		
<i>Submit to board</i>			✓	
Wrap Up			✓	✓
Preparing the presentation of the thesis			✓	✓
Thesis defense				✓

Reference

-  Asim, M. et al. (2019) **Two Stream Deep Network for Document Image Classification**. International Conference on Document Analysis and Recognition (ICDAR), 2019, pp. 1410-1416, doi: 10.1109/ICDAR.2019.00227.
-  Audebert N. et al. (2020) **Multimodal Deep Networks for Text and Image-Based Document Classification**. In: Cellier P., Driessens K. (eds) Machine Learning and Knowledge Discovery in Databases (ECML/PKDD) Communications in Computer and Information Science, Springer, March, 2020, vol 1167: 427-443.
-  Braz, F. A. et al. (2021). **Leveraging effectiveness and efficiency in Page Stream Deep Segmentation**, Engineering Applications of Artificial Intelligence, Volume 105, 2021, ISSN 0952-197, <https://doi.org/10.1016/j.engappai.2021.104394>.
-  Devlin, J., Chang, M., Lee, K. and Toutanova, K. (2019) **BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding**. NAACL-HLT (1) 2019: 4171-4186
-  Howard, Jeremy and Ruder, Sebastian. (2018). **Universal Language Model Fine-tuning for Text Classification**. 328-339. 10.18653/v1/P18-1031.
-  Merity, S., Keskar, N. and Socher, R. (2018). **Regularizing and Optimizing LSTM Language Models**. International Conference on Learning Representations, 2018: 1-13.
-  Wiedemann, G. and Heyer, G. (2021). **Multi-modal page stream segmentation with convolutional neural networks**. Language Resources and Evaluation. 55. 10.1007/s10579-019-09476-2.

Thanks!
Questions?
Suggestion?



Patricia Medyna

patriciamedyna@aluno.unb.br