



D'ici, on voit + loin !

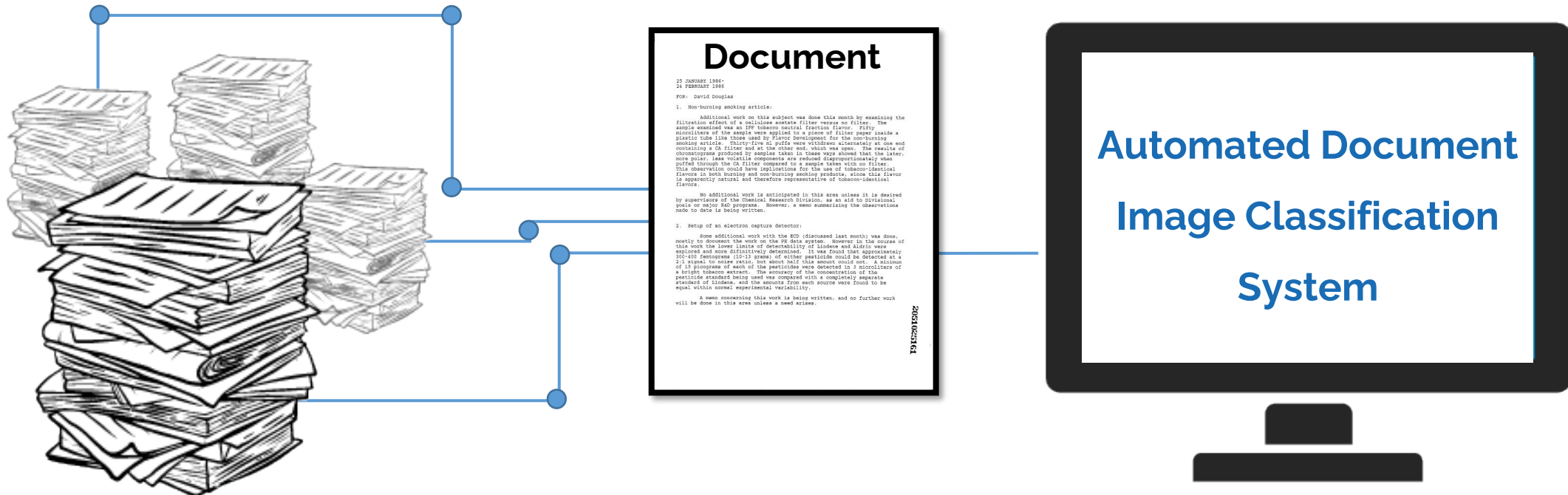
EAML: Ensemble Self-Attention-based Mutual Learning Network for Document Image Classification

Souhail Bakkali, Mickaël Coustaty, Zuheng Ming, Marçal Rusiñol, Oriol Ramos Terrades

14/10/2022



Why do we need Multimodal Document Image Classification ?



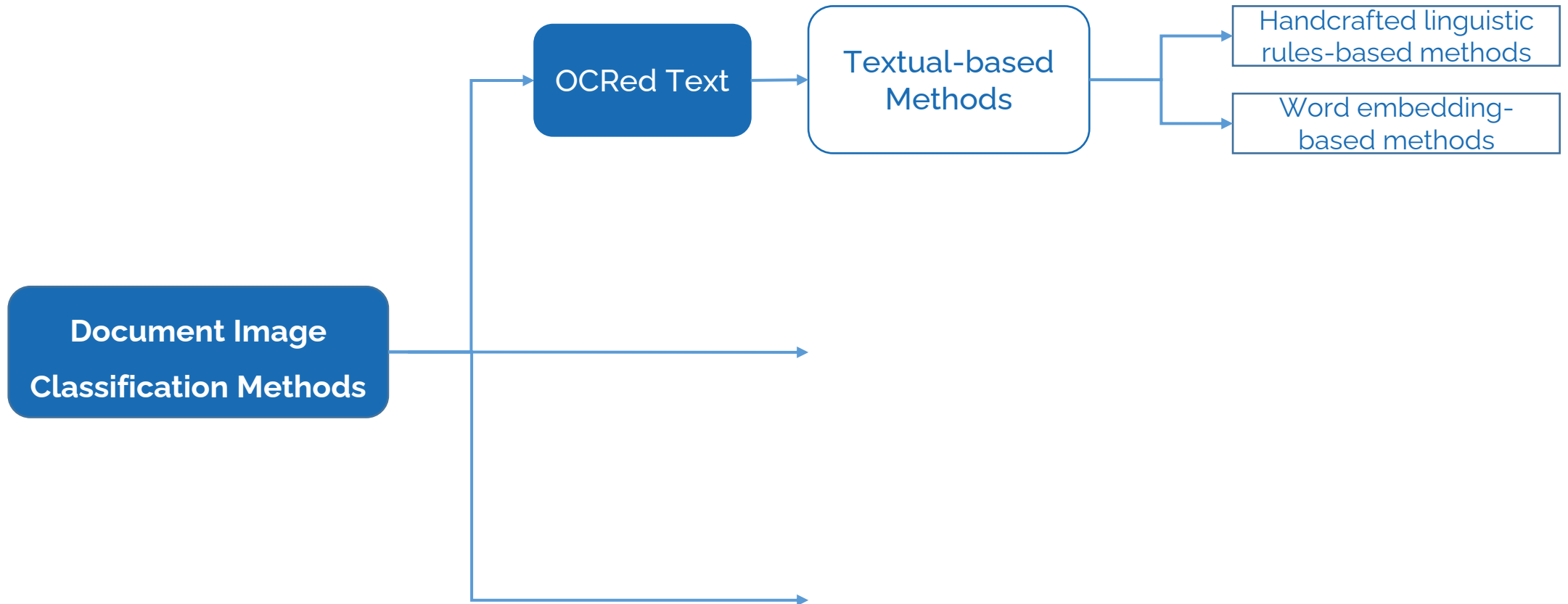
The Existing Document Image classification Methods

Document Image
Classification Methods

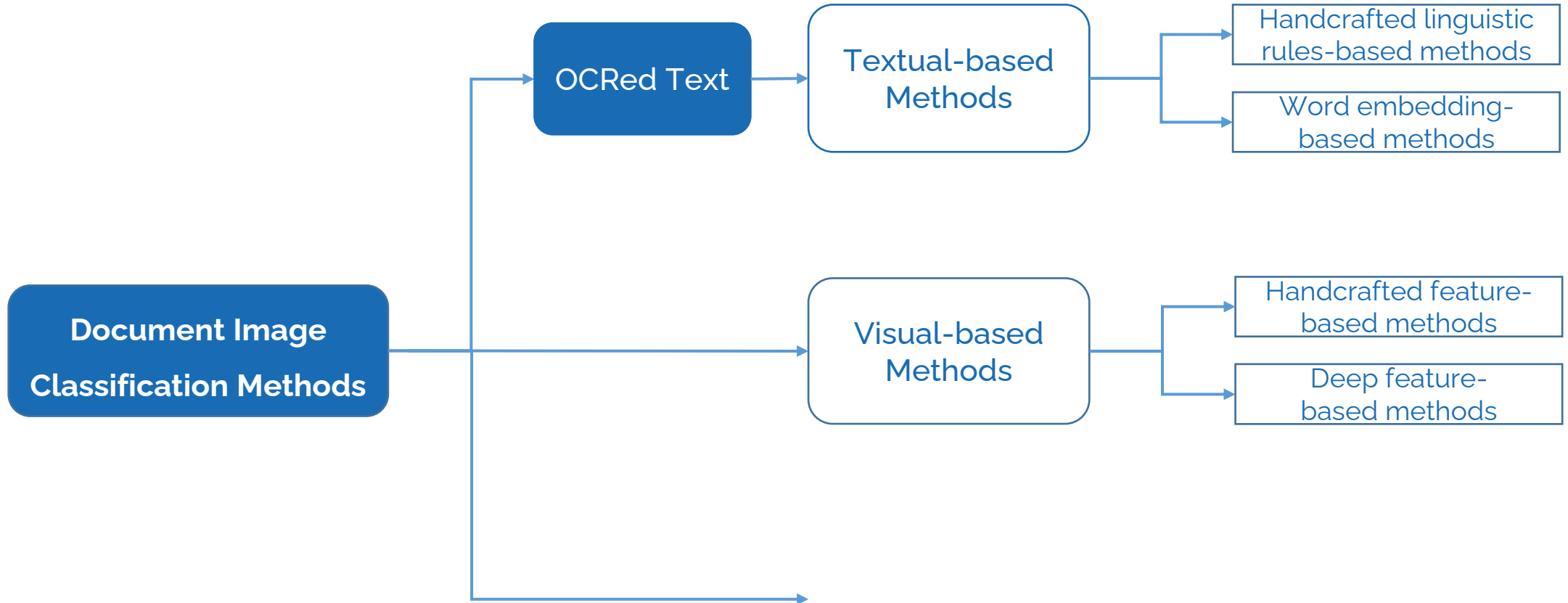


```
graph LR; A[Document Image Classification Methods] --> B[ ]; A --> C[ ]; A --> D[ ]
```

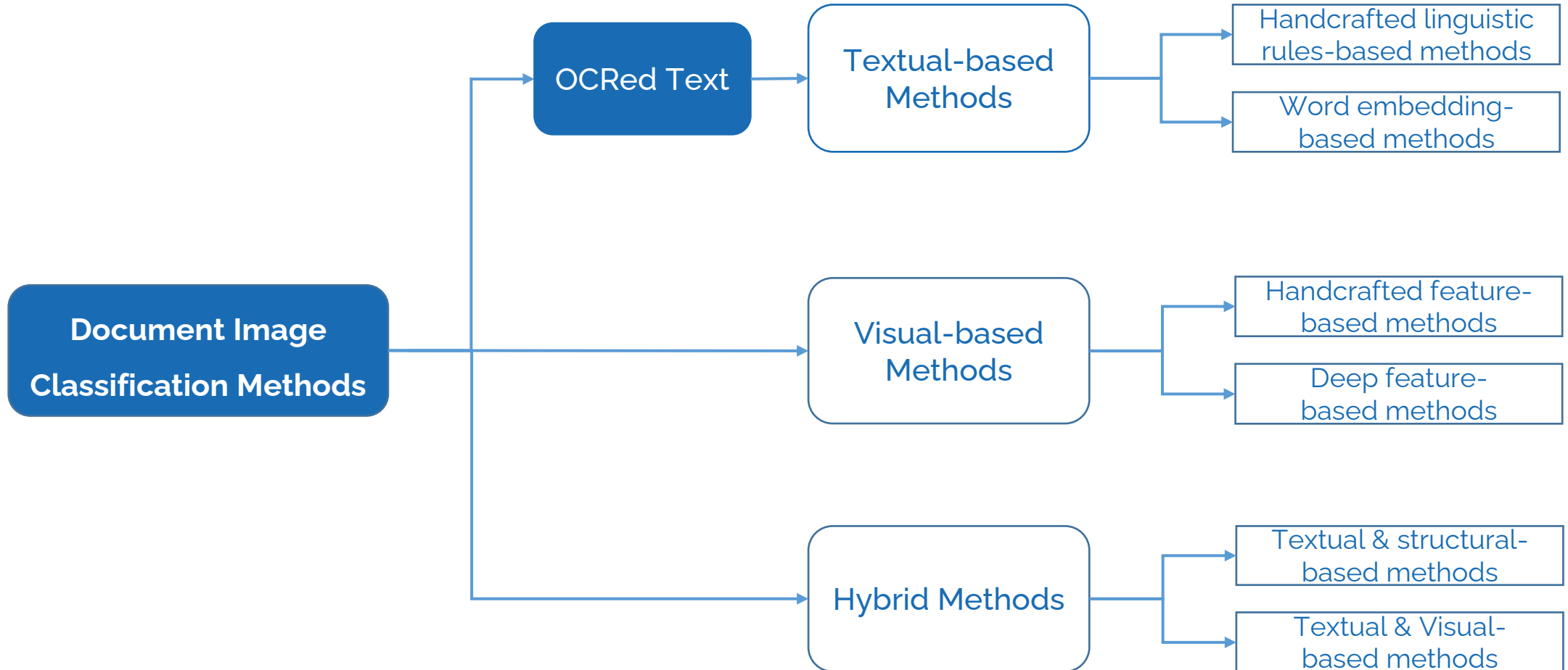
The Existing Document Image classification Methods



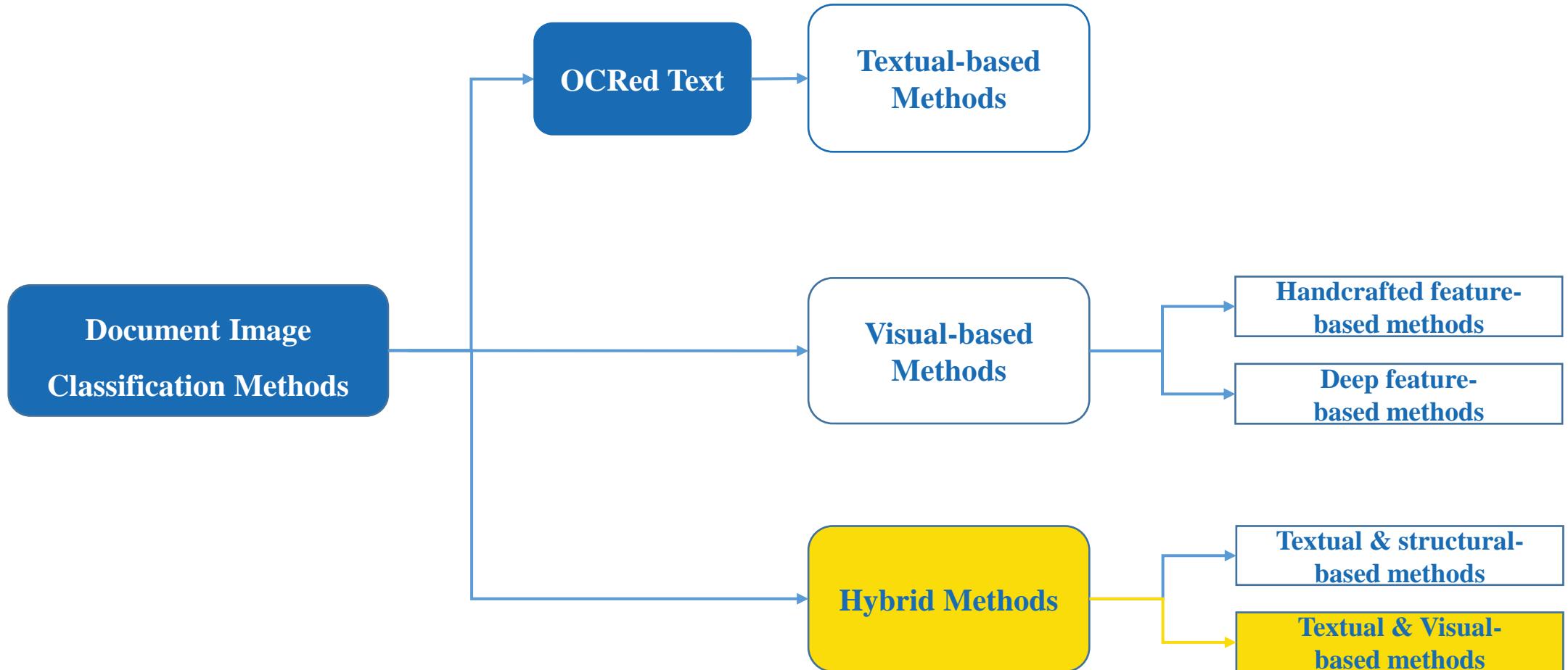
The Existing Document Image classification Methods

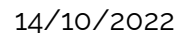


The Existing Document Image classification Methods



The Existing Document Image classification Methods





Our Proposed EAML Network

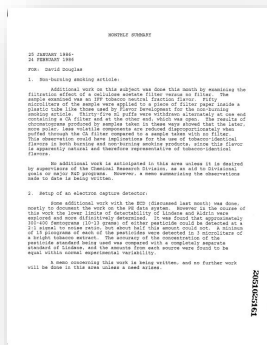


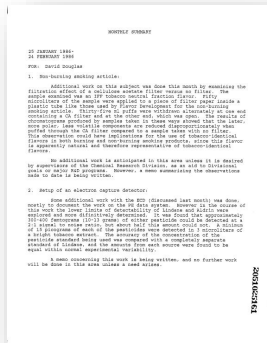
Image Branch
(with attention module)

OCRed
Text

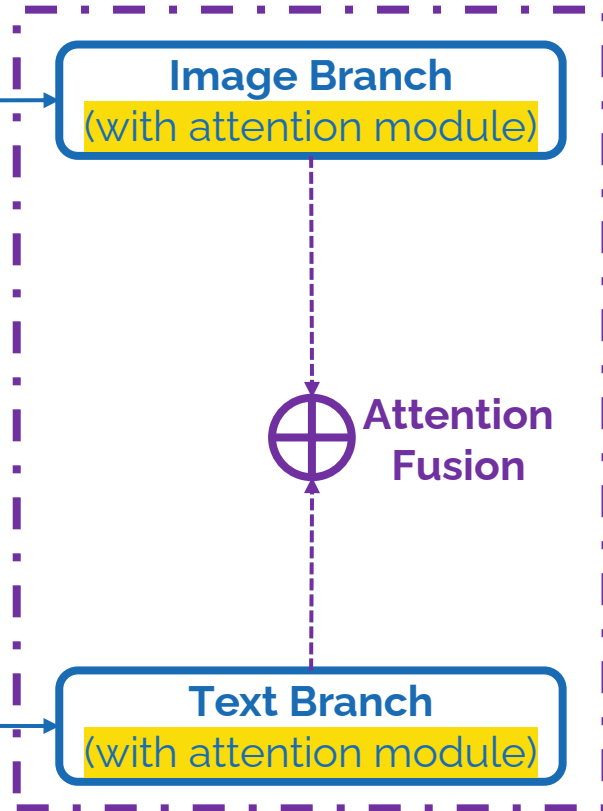


Text Branch
(with attention module)

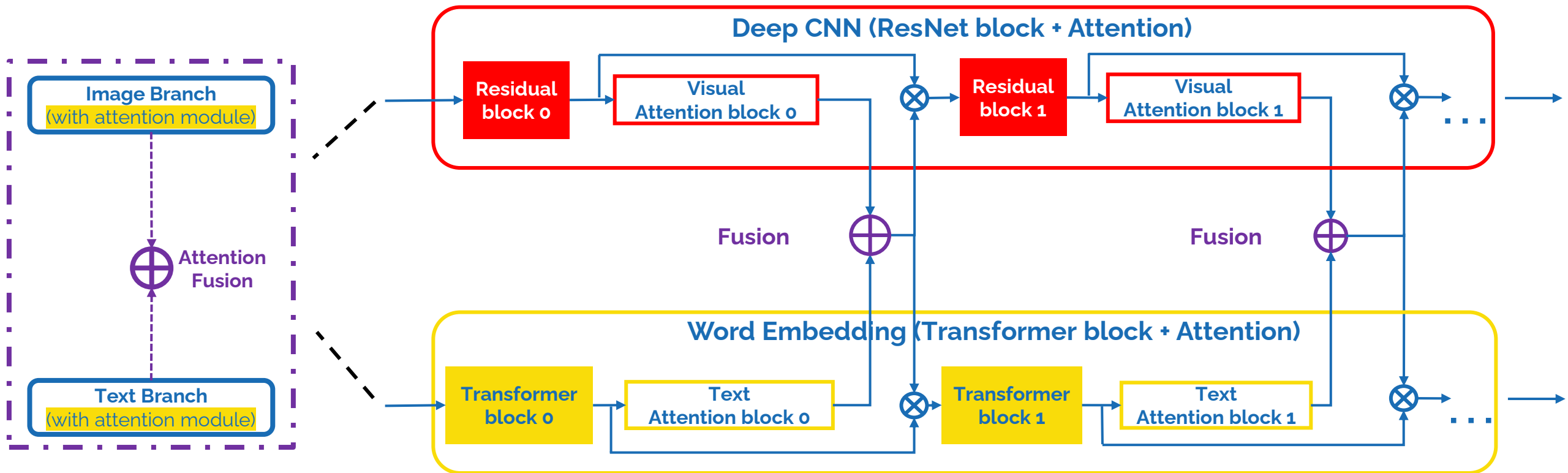
Our Proposed EAML Network

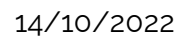


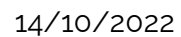
OCRed
Text



Self-Attention-based Fusion Module

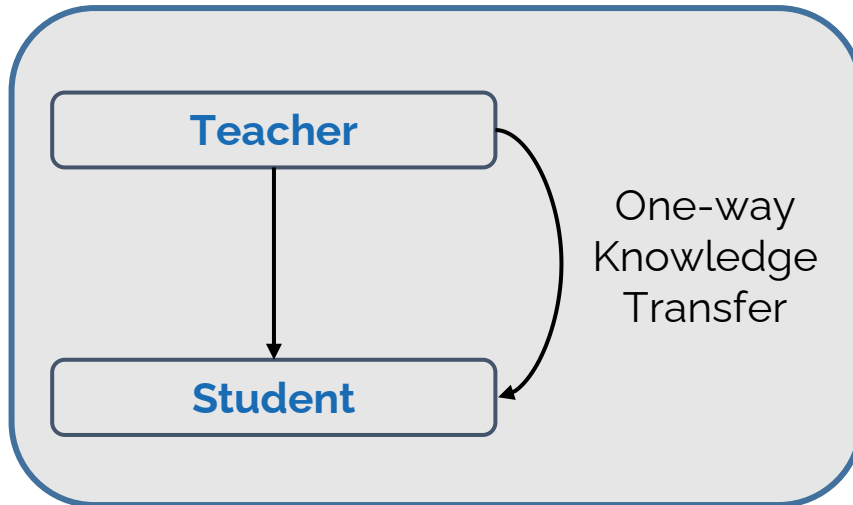






Conventional Mutual Learning in Computer Vision

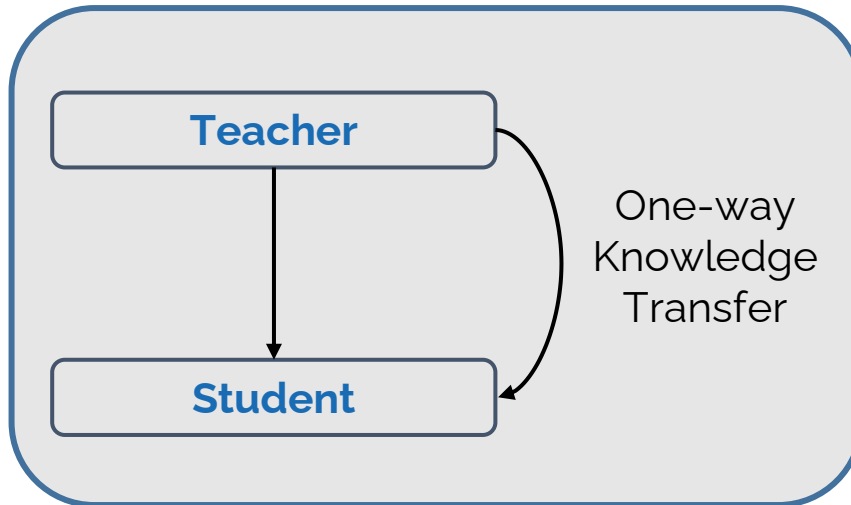
Knowledge Distillation-based approach



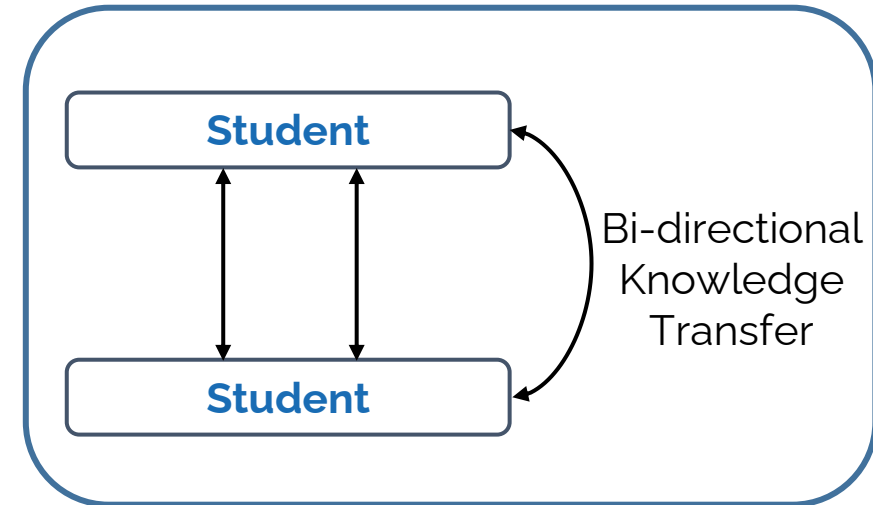
Hinton, Geoffrey E., Oriol Vinyals and J. Dean. "Distilling the Knowledge in a Neural Network." ArXiv abs/1503.02531 (2015): n. pag.

Conventional Mutual Learning in Computer Vision

Knowledge Distillation-based approach



Conventional Mutual Learning-based approach

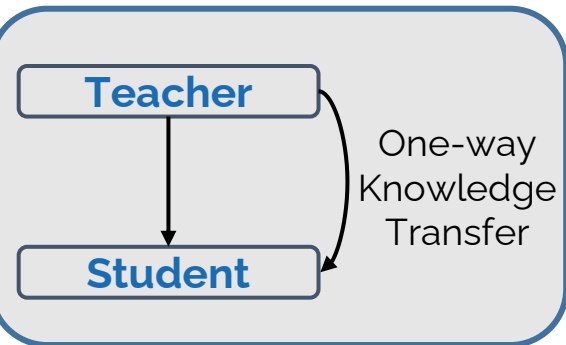


Hinton, Geoffrey E., Oriol Vinyals and J. Dean. "Distilling the Knowledge in a Neural Network." ArXiv abs/1503.02531 (2015): n. pag.

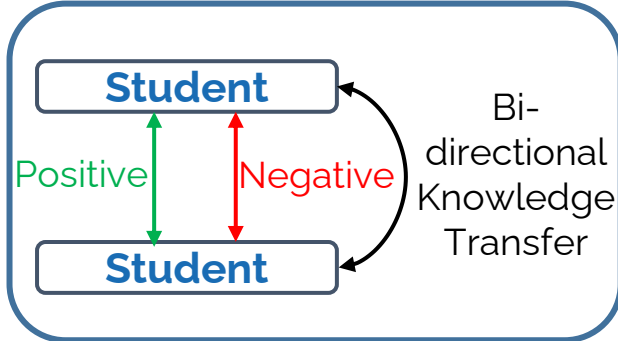
Zhang, Y., T. Xiang, Timothy M. Hospedales And Huchuan Lu. "Deep Mutual Learning." 2018 IEEE/CVF Conference On Computer Vision And Pattern Recognition (2018): 4320-4328.

Conventional Mutual Learning in Computer Vision: Challenges

Knowledge Distillation- based approach



Conventional Mutual Learning-based approach



Encourages collaborative learning between modalities.



Enables to minimize the difference in class probabilities produced by the two modalities



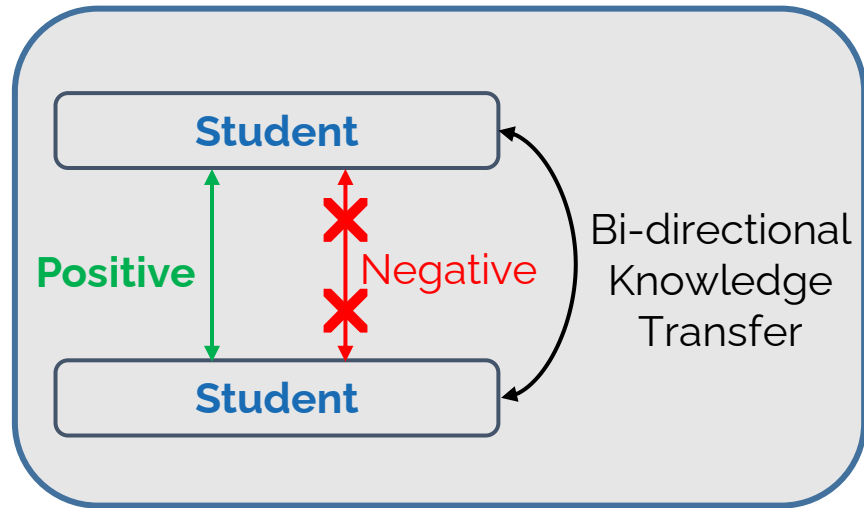
Both modalities can learn both the positive & negative information from one another



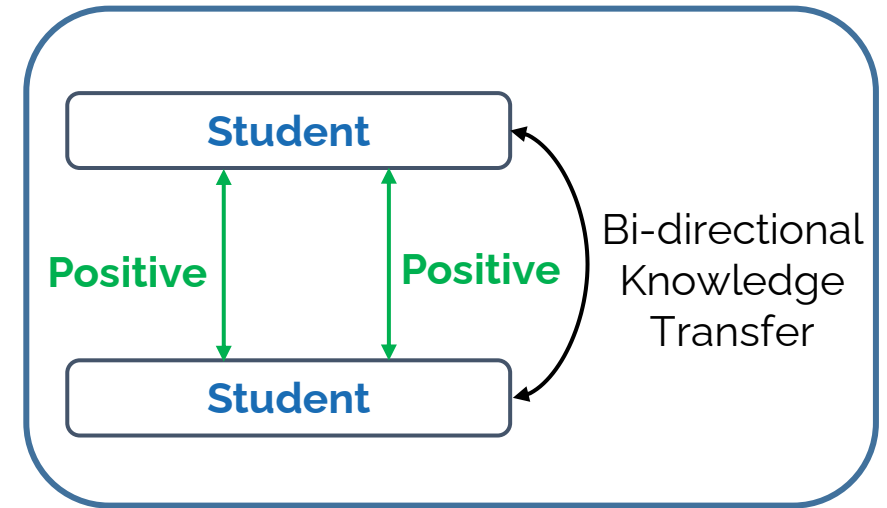
The introduction of the negative knowledge from one modality to another harms and weakens the ongoing training.

Our Proposed Positive Mutual Learning Strategy

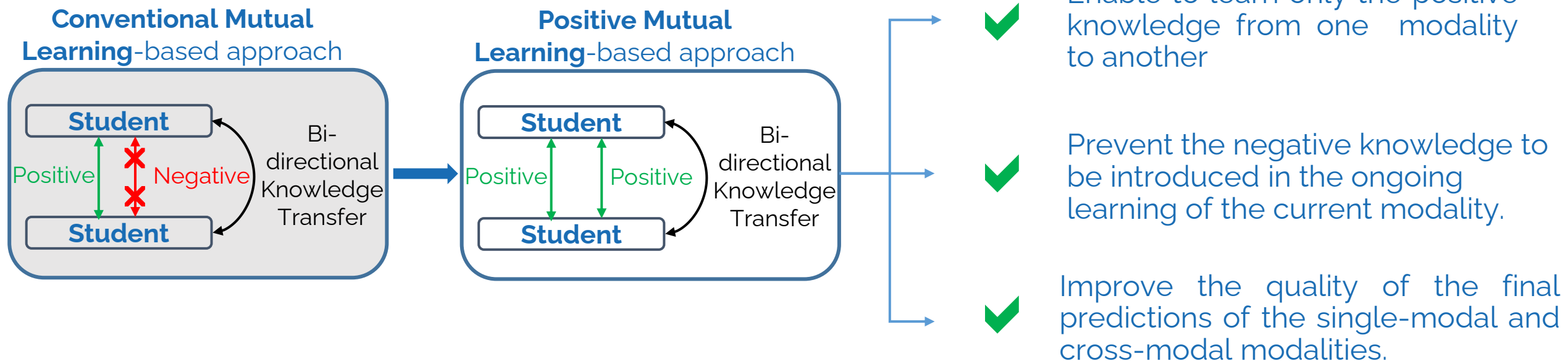
Conventional Mutual Learning-based approach



Positive Mutual Learning-based approach (Ours)



Our Proposed Positive Mutual Learning Strategy



Experimental Setup & Ablation Study

+ Experimental Setup

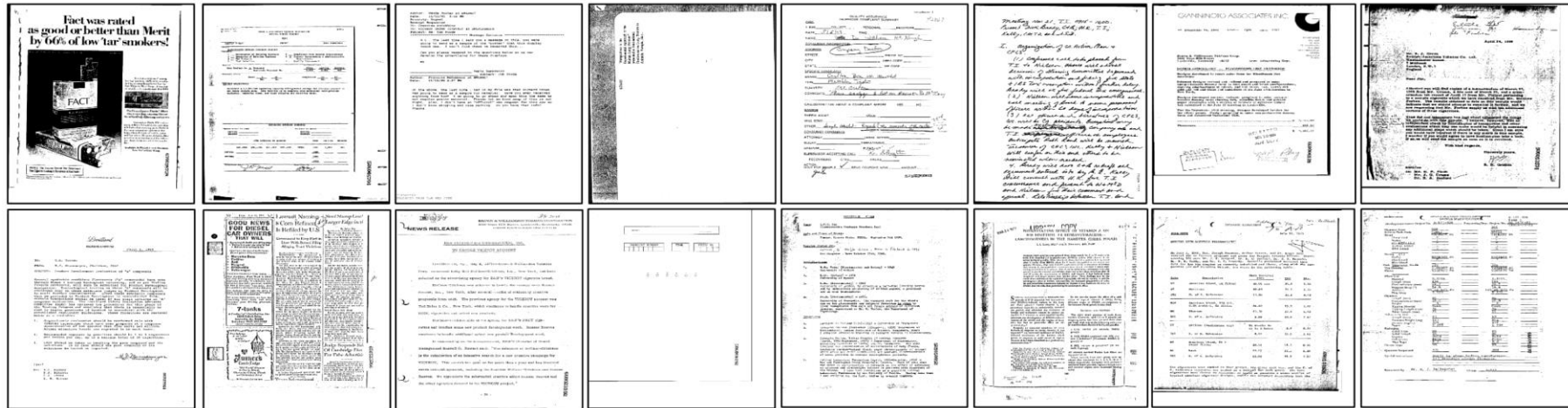
> Datasets

- RVLCDIP
- Tobacco-3482

+ Ablation Study

- > Intra-Dataset Evaluation
- > Inter-Dataset Evaluation
- > Intra-Dataset & Inter-Dataset Evaluation comparison

Experimental Setup: **Datasets**

**RVLCDIP**

- **400,000** gray-scale document images
- Dataset split:
 - Training set: 320,000 images.
 - Validation set: 40,000 images.
 - Test set: 40,000 images.
- Samples of different document classes in the **RVL-CDIP** dataset. From left to right: Advertisement, Budget, Email, File folder, Form, Handwritten, Invoice, Letter, Memo, News article, Presentation, Questionnaire, Resume, Scientific publication, Scientific report, Specification.

Experimental Setup: Datasets

Tobacco-3482



- **3482** gray-scale document images.
- Samples of different document classes in the **Tobacco-3482** dataset. The 10 categories are: *ADVE*, *Email*, *Form*, *Letter*, *Memo*, *News*, *Notes*, *Report*, *Resume*, *Scientific*.
- Dataset Split:
 - 80% Training set
 - 10% Validation set
 - 10% Test set.

Ablation Study : Intra-Dataset Evaluation

Overall Classification Accuracy (%) on the proposed approaches

RVLCDIP	Method	Image Modality	Text Modality	Fusion Modality
	Independent Learning	85.04	84.96	94.44

Ablation Study : Intra-Dataset Evaluation

Overall Classification Accuracy (%) on the proposed approaches

RVLCDIP	Method	Image Modality	Text Modality	Fusion Modality
	Independent Learning	85.04	84.96	94.44
	Conventional Mutual Learning	88.87	80.89	90.06

Ablation Study : Intra-Dataset Evaluation

Overall Classification Accuracy (%) on the proposed approaches

RVLCDIP	Method	Image Modality	Text Modality	Fusion Modality
	Independent Learning	85.04	84.96	94.44
	Conventional Mutual Learning	88.87	80.89	90.06
	Positive Mutual Learning (Ours)	90.81	88.80	96.28

Ablation Study : Intra-Dataset Evaluation

Overall Classification Accuracy (%) on the proposed approaches

RVLCDIP	Method	Image Modality	Text Modality	Fusion Modality
	Independent Learning	85.04	84.96	94.44
	Conventional Mutual Learning	88.87	80.89	90.06
	Positive Mutual Learning (Ours)	90.81	88.80	96.28
	EAML	97.67	97.63	97.70

Ablation Study : Intra-Dataset Evaluation

Overall Classification Accuracy (%) on the proposed approaches

Tobacco-3482	Method	Image Modality	Text Modality	Fusion Modality
	Independent Learning	96.17	96.02	96.95
	Conventional Mutual Learning	93.69	88.82	94.84
	Positive Mutual Learning (Ours)	97.70	96.27	98.28
	EAML	97.99	96.27	98.57

Ablation Study : Intra-Dataset Evalutation

Overall Classification Accuracy (%) on the proposed approaches

RVLCDIP

Method	Model	Accuracy(%)
Image	Nicolas <i>et al.</i> [7]	89.1
Text		74.6
Multi-modal		90.6
Image	Dauphinee <i>et al.</i> [17]	90.24
Text		82.23
Multi-Modal		93.07
Image	Cross-Modal [11]	91.45
Text		84.96
Multi-modal		97.05
Image	EAML _{Tr-KLD_{Reg}} (Ours)	97.67
Text		97.63
Multi-Modal		97.70
Baselines	Harley <i>et al.</i> [25]	89.80
	Csurka <i>et al.</i> [15]	90.70
	Tensmeyer <i>et al.</i> [56]	90.94
	Azfal <i>et al.</i> [2]	90.97
	Das <i>et al.</i> [16]	91.11
	Das <i>et al.</i> [16]	92.21
	Ferrando <i>et al.</i> [21]	92.31
	Xu <i>et al.</i> [61]	94.42
	Xu <i>et al.</i> [62]	95.64

Tobacco-3482

Method	Model	Accuracy(%)
Image	Nicolas <i>et al.</i> [7]	84.5
Text		73.8
Multi-modal		87.8
Image	Asim <i>et al.</i> [6]	93.2
Text		87.1
Multi-modal		95.8
Image	Ferrando <i>et al.</i> [21]	94.04
Text		-
Multi-modal		94.90
Image	Cross-Modal [10]	96.25
Text		97.18
Multi-modal		99.71
Image	EAML _{Tr-KLD_{Reg}} (Ours)	97.99
Text		96.27
Multi-modal		98.57
Baselines	Kumar <i>et al.</i> [34]	43.8
	Kang <i>et al.</i> [30]	65.37
	Afzal <i>et al.</i> [1]	76.6
	Harley <i>et al.</i> [25]	79.9
	Noce <i>et al.</i> [43]	79.8

Ablation Study: Inter-Dataset Evaluation Protocol

Inter-Dataset	Train	Test
	RVLCDIP	Tobacco-3482
	Tobacco-3482	RVLCDIP

Ablation Study :

Intra-Dataset & Inter-Dataset Evaluation Comparison

Train on Tobacco-3482 ▶ ▶ Test on Tobacco-3482

Method	Image Modality	Text Modality	Fusion Modality
Positive Mutual Learning (Ours)	97.70	96.27	98.28
EAML	97.99	96.27	98.57

Train on RVLCDIP ▶ ▶ Test on Tobacco-3482

Method	Image Modality	Text Modality	Fusion Modality
Positive Mutual Learning (Ours)	84.82	83.73	86.68
EAML	87.29	87.23	87.63

Ablation Study :

Intra-Dataset & Inter-Dataset Evaluation Comparison

● Train on RVLCDIP ▶ ▶ Test on RVLCDIP

Method	Image Modality	Text Modality	Fusion Modality
EAML	97.67	97.63	97.70

● Train on Tobacco-3482 ▶ ▶ Test on RVLCDIP

Method	Image Modality	Text Modality	Fusion Modality
EAML	78.89	79.06	86.68



D'ici, on voit + loin !

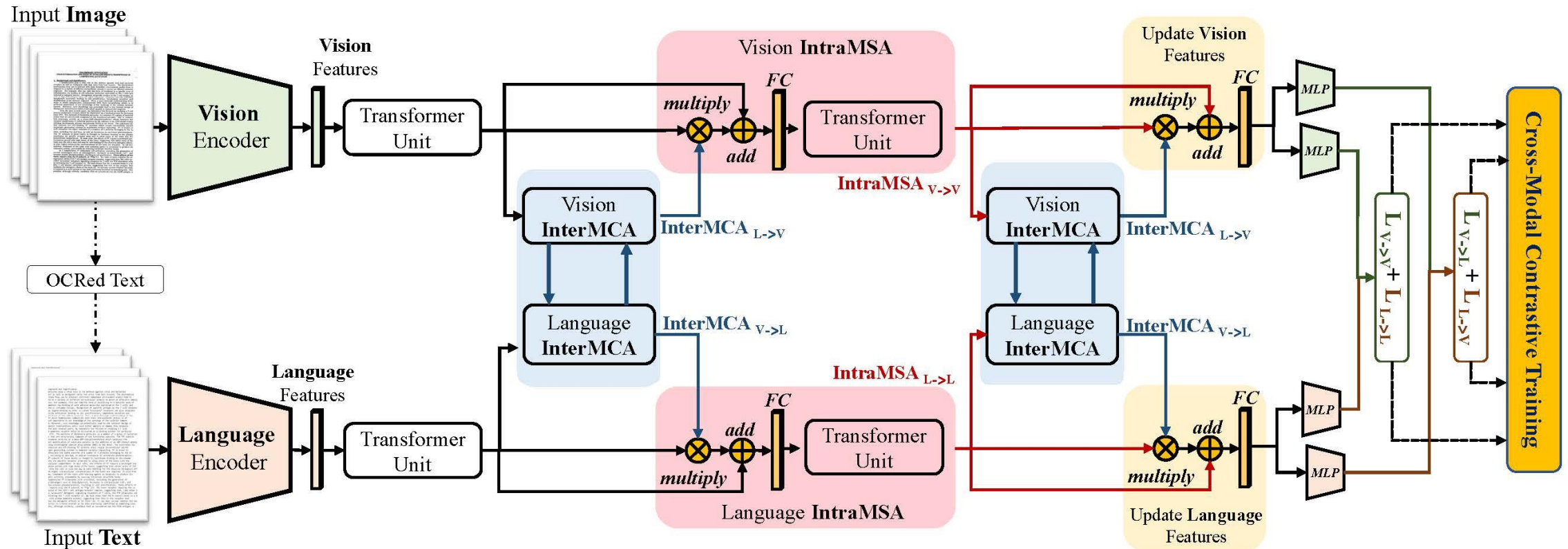


univ-larochelle.fr

Current Work: Cross-Modal Learning (i.e. pre-training) to perform Document Classification (as a downstream task)

- + We approach the document classification problem by learning cross-modal representations through language and vision cues, considering **intra- and inter-modality relationships**.
- + Instead of merging features from different modalities into a common representation space, we intend to exploit high-level interactions and learn semantic information from effective attention flows **within and across** modalities.
- + The learning objective that we propose aims to enforce the compactness of intra-class representations while separating inter-class features by contrasting positive and negative sample pairs within and across each modality.

Current Work: Cross-Modal Learning (i.e. pre-training) to perform Document Classification (as a downstream task)



- + Overview of the proposed cross-modal contrastive learning method. The network is composed of InterMCA and IntraMSA modules with flexible attention mechanisms to learn cross-modal representations in a cross-modal contrastive learning fashion.

Current Work: Cross-Modal Alignment

- + We introduce the InterMCA and IntraMSA attention modules that capture intrinsic patterns by modeling the inter-modality and intra-modality relationships for image regions and texts.
- + **Inter-Modal Alignment:**
 - > The inter-modality cross-attention module **InterMCA** aims to enhance the cross-modal features by embracing cross-modal interactions across image regions and texts. This module aims to transfer the salient information from one modality to another.

$$\mathbf{InterMCA}_{\mathbf{L} \rightarrow \mathbf{V}}(\mathbf{V}^l) = \text{softmax} \left(\frac{\mathcal{Q}_{\mathbf{V}^l} \mathcal{K}_{\mathbf{L}^l}^\top}{\sqrt{d_k}} \right) \mathcal{V}_{\mathbf{L}^l} \quad (1)$$

$$\mathbf{InterMCA}_{\mathbf{V} \rightarrow \mathbf{L}}(\mathbf{L}^l) = \text{softmax} \left(\frac{\mathcal{Q}_{\mathbf{L}^l} \mathcal{K}_{\mathbf{V}^l}^\top}{\sqrt{d_k}} \right) \mathcal{V}_{\mathbf{V}^l} \quad (2)$$

Current Work: Cross-Modal Alignment

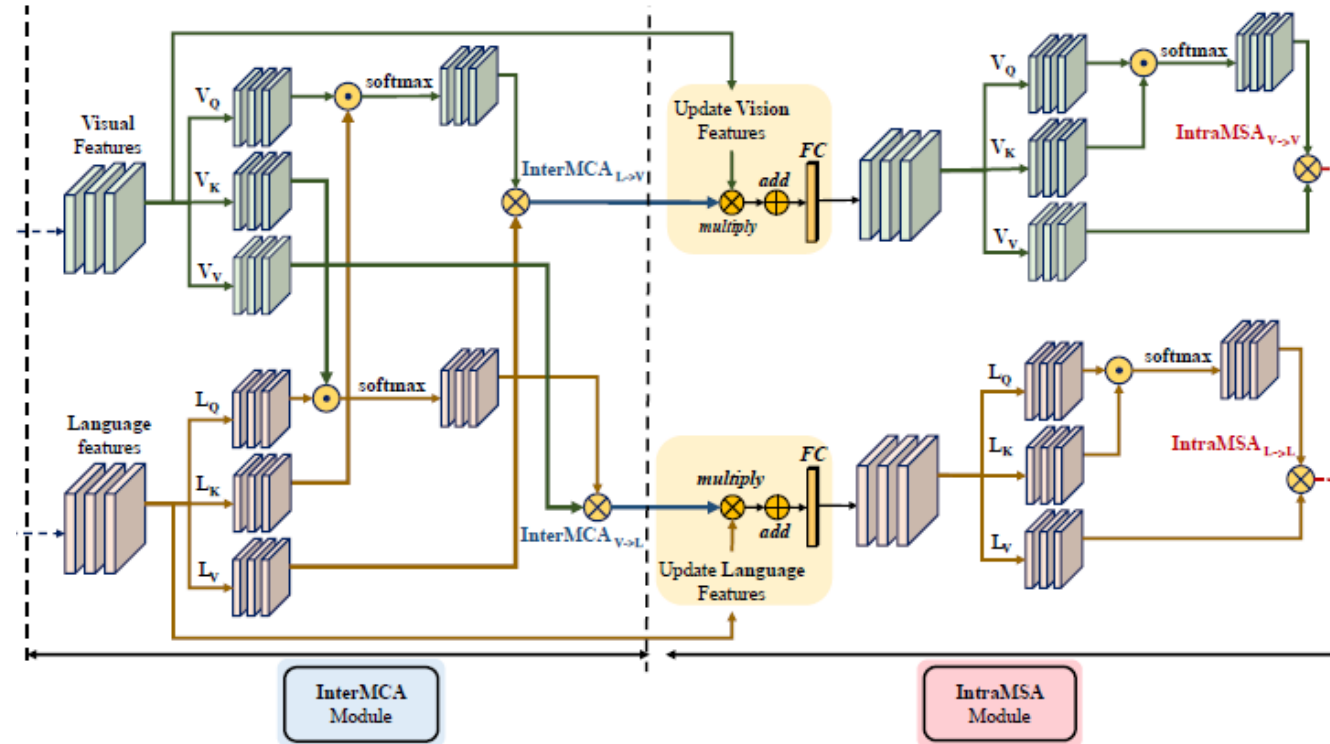
+ Intra-Modal Alignment:

- > The intra-modal self-attention module **IntraMSA** aims to update the vision and language information and to capture inner-modality attention weights. For each modality, the information is updated according to a feature fusion scheme.

$$\mathbf{IntraMSA}_{V \rightarrow V} = \text{softmax} \left(\frac{Q_{\hat{V}^l} K_{\hat{V}^l}^\top}{\sqrt{d_k}} \right) V_{\hat{V}^l} \quad (9)$$

$$\mathbf{IntraMSA}_{L \rightarrow L} = \text{softmax} \left(\frac{Q_{\hat{L}^l} K_{\hat{L}^l}^\top}{\sqrt{d_k}} \right) V_{\hat{L}^l} \quad (10)$$

Current Work: Cross-Modal Alignment

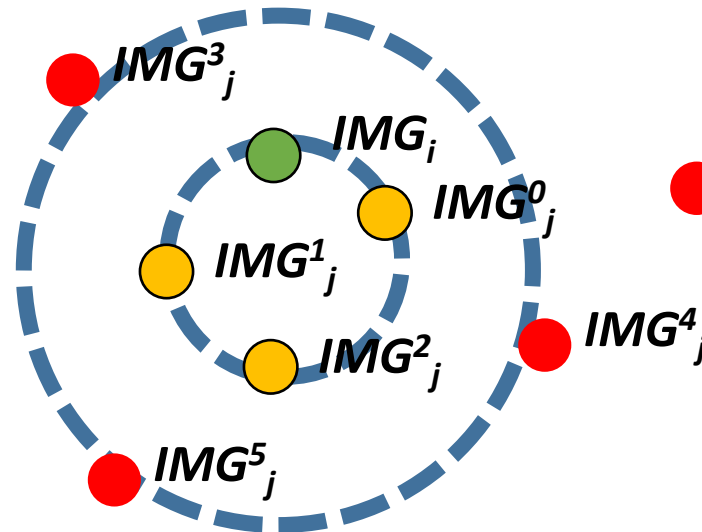


- + Illustration of the **InterMCA** and **IntraMSA** attention modules. The visual and textual features are transformed into query, key, and value vectors. They are jointly leveraged and are further fused to transfer attention flows between modalities to update the original features.

Current Work: Cross-Modal Alignment

+ Intra-Modality Learning Objective

● Positive Pairs for each Image are (IMG_i, IMG_j)

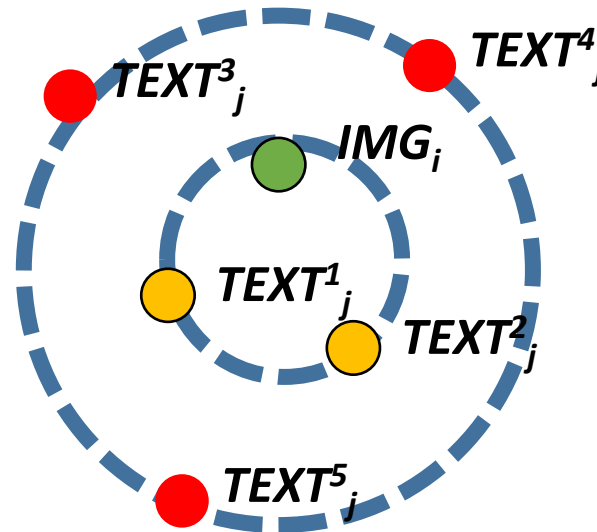


● Negative Pairs for each Image are (IMG_i, IMG_k)

Current Work: Cross-Modal Alignment

+ Inter-Modality Learning Objective

● Positive Pairs for each Image are $(IMG_i, TEXT_j)$



● Negative Pairs for each Image are $(IMG_i, TEXT_k)$

Pre-training setting	IntraMSA	InterMCA	#Param	Acc.(%)
<i>-w/o lang modality</i>				
			198M	85.71
	✓		201M	86.66
		✓	209M	87.20
	✓	✓	217M	90.94
<i>-w/o visn modality</i>				
			198M	86.01
	✓		201M	86.31
		✓	209M	87.50
	✓	✓	217M	90.62

Table 1: Ablation study on CMDoC on cross-modality attention components, pre-trained on Tobacco dataset

Model	Accuracy (%)		
	T → R	R → T	R → N
<i>w/o language modality</i>			
- EAML [Bakkali <i>et al.</i> , 2021]	78.89	84.82	-
- CMDoC	79.04	89.73	99.99
<i>w/o vision modality</i>			
- EAML [Bakkali <i>et al.</i> , 2021]	79.06	83.72	-
- CMDoC	81.96	89.88	99.99

Table 3: Cross-dataset test on datasets with different size and document types. T,R, and N denote Tobacco, RVL-CDIP, and NIST datasets. T → R denotes fine-tune on Tobacco, and test on RVL-CDIP.

Current Work: Results

Method	Accuracy(%)
<i>vision only methods</i>	
VGG-16 [Afzal <i>et al.</i> , 2017]	90.31
AlexNet [Tensmeyer and Martinez, 2017]	90.94
Ensemble [Das <i>et al.</i> , 2018]	92.21
<i>language only methods</i>	
BERT _{Base} [Devlin <i>et al.</i> , 2019]	86.10
RoBERTa _{Base} [Liu <i>et al.</i> , 2019]	90.94
LayoutLM _{Base} [Xu <i>et al.</i> , 2020]	90.11
<i>vision+language methods</i>	
w/o language	
- Multimodal [Audebert <i>et al.</i> , 2019]	89.1
- Ensemble [Dauphinee <i>et al.</i> , 2019]	91.45
- EAML [Bakkali <i>et al.</i> , 2021]	90.81
w/o vision	
- Multimodal [Audebert <i>et al.</i> , 2019]	74.6
- Ensemble [Dauphinee <i>et al.</i> , 2019]	82.23
- EAML [Bakkali <i>et al.</i> , 2021]	88.80
CMDoC (V+L) w/o language	92.64
CMDoC (V+L) w/o vision	91.37

CMDoC (V+L) w/o language	92.64
CMDoC (V+L) w/o vision	91.37
<i>vision+language+layout methods</i>	
TILT _{Base} [Powalski <i>et al.</i> , 2021]	93.50
SelfDoc [Li <i>et al.</i> , 2021]	93.81
LayoutLM _{Base} [Xu <i>et al.</i> , 2020]	94.42
LayoutLMv2 _{Base} [Xu <i>et al.</i> , 2022]	95.25
DocFormer _{Base} [Appalaraju <i>et al.</i> , 2021]	96.17

Table 4: Top-1 accuracy (%) comparison results of different document classification methods evaluated on the of RVL-CDIP dataset. V+L denotes vision+language modalities

Future Work

- + Essentially, we are exploring Self-Supervised Pre-training techniques to develop a more general and domain-agnostic multi-modal embedding network to be fine-tuned on downstream applications such as: few-shot document classification, document retrieval, domain generalization, etc.

+ Accepted Papers:

- > S. Bakkali, Z. Ming, M. Coustaty and M. Rusiñol, "Cross-Modal Deep Networks For Document Image Classification," 2020 IEEE International Conference on Image Processing (ICIP), 2020, pp. 2556-2560, doi: 10.1109/ICIP40778.2020.9191268.
- > S. Bakkali, Z. Ming, M. Coustaty and M. Rusiñol, "Visual and Textual Deep Feature Fusion for Document Image Classification," *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2020, pp. 2394-2403, doi: 10.1109/CVPRW50498.2020.00289.
- > Bakkali, S., Ming, Z., Coustaty, M. et al. EAML: ensemble self-attention-based mutual learning network for document image classification. IJDAR 24, 251–268 (2021). <https://doi.org/10.1007/s10032-021-00378-0>

+ In Progress:

- > S. Bakkali, Z. Ming, M. Coustaty, M. Rusiñol, and OR. Terrades, "CMDDoC: Cross-Modal Learning for Document Classification". (Under Review: Pattern Recognition Journal)



D'ici, on voit + loin !



univ-larochelle.fr



D'ici, on voit + loin !

EAML: Ensemble Self-Attention-based Mutual Learning Network for Document Image Classification

Souhail Bakkali, Mickaël Coustaty, Zuheng Ming, Marçal Rusiñol, Oriol Ramos Terrades

14/10/2022

