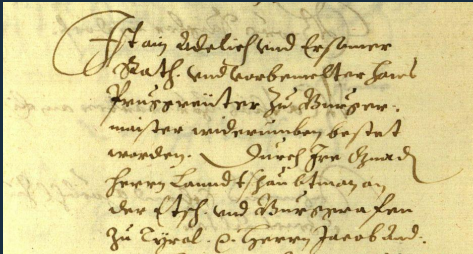
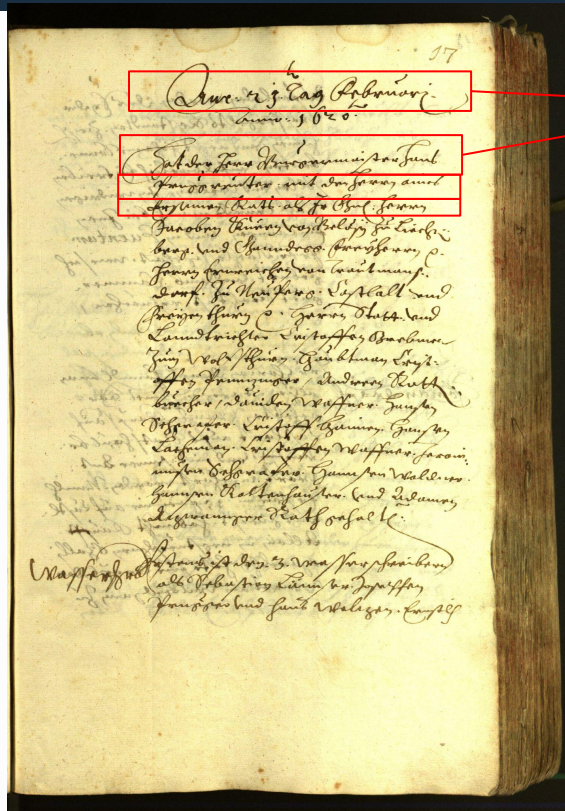


# Apprentissage semi-supervisé multi-hypothèses pour la reconnaissance d'écriture manuscrite

Alexandre Chapin, Yann Soullard, Bertrand Couasnon, Aurélie Lemaître



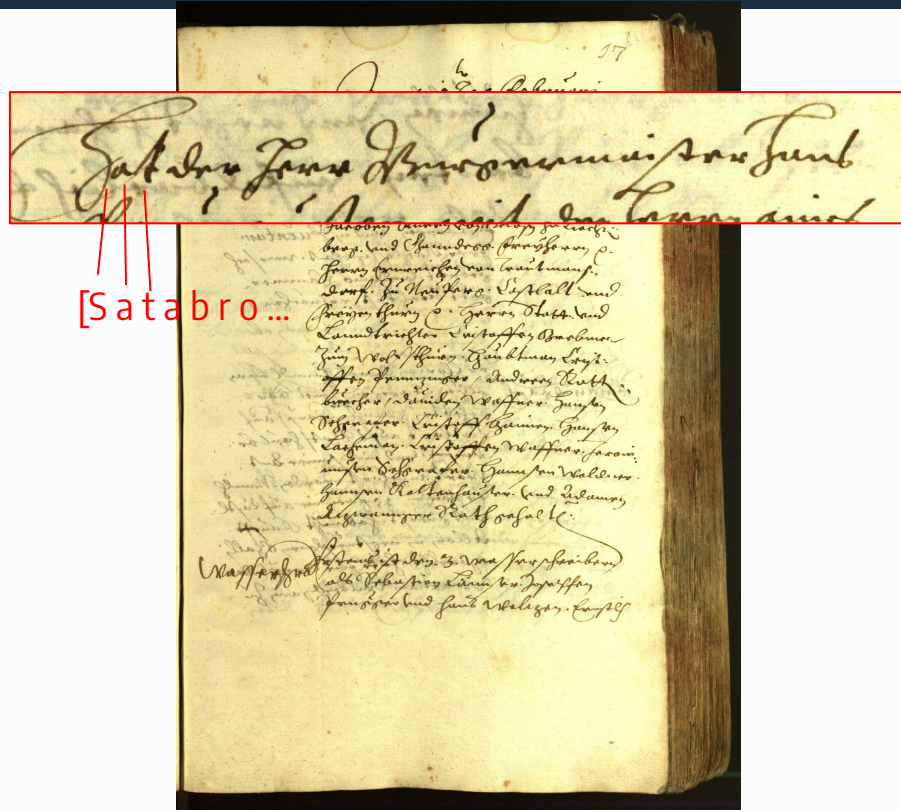
# Contexte : reconnaissance d'écriture manuscrite



Identification des  
zones de texte

Image extraite de la base de données READ18

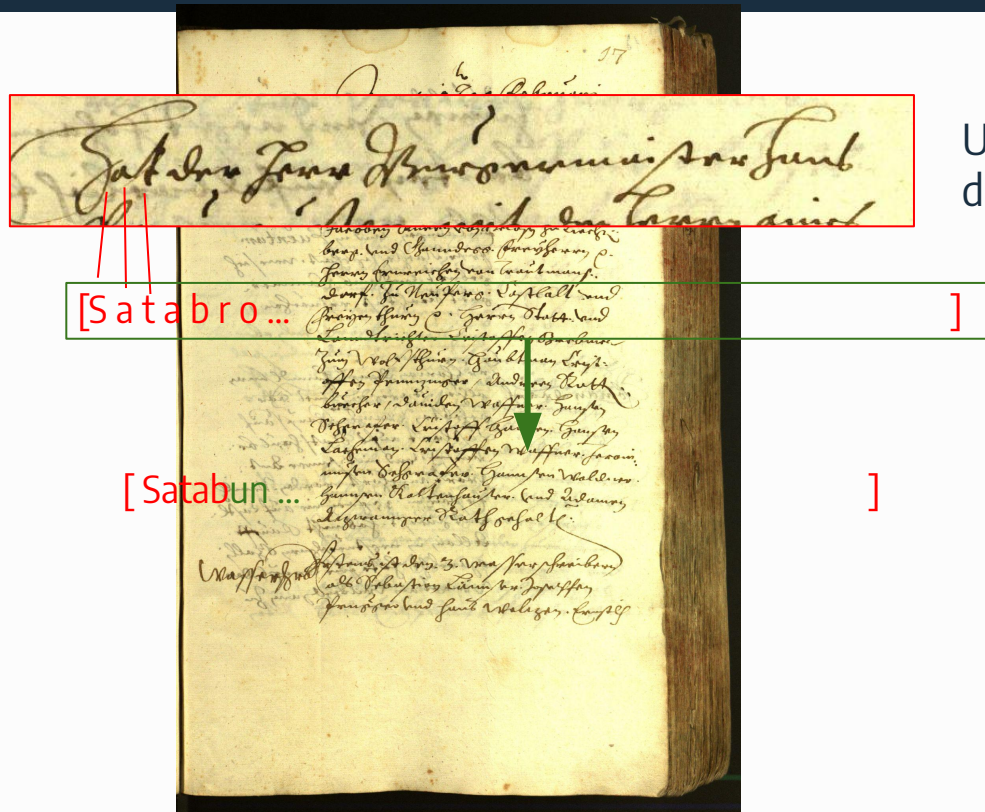
# Contexte : reconnaissance d'écriture manuscrite



Reconnaissance de texte

Image extraite de la base de données READ18

# Contexte : reconnaissance d'écriture manuscrite

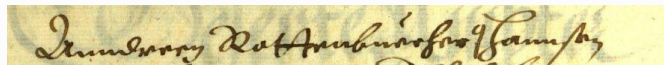


Utilisation d'un modèle de langage :

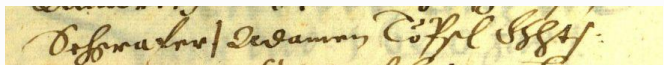
- Pour corriger les erreurs

# Le problème du manque de données

Lignes de texte avec  
**transcription** associées  
**(annotations) :**



*'Anndreen Rottenbuecher. Hannsen'*



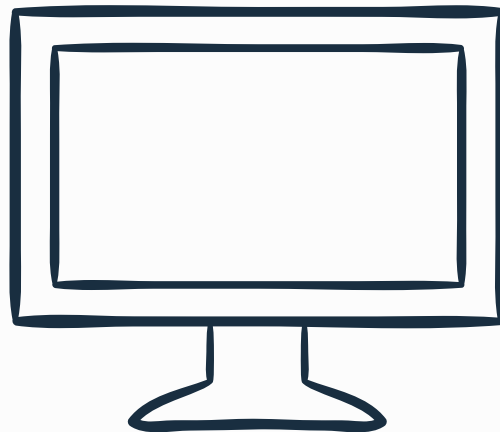
*'Schgrafer/ Adamen TōPsl Rhts-'*

Peu de données annotées car difficile à  
annoter (même pour un expert)

Ont besoin de ...



Modèles d'**apprentissage  
profond**



# Notre contribution

- **Combiner des méthodes** d'apprentissage profond
- Apprendre avec **peu de données**
- Tirer profit d'**exemples non étiquetés**

01

Etat de l'art

# Modèles classiques : CRNN

## Convolutionnal Recurrent Neural Network

### Principe

- **Couches de convolution** en entrée : extraction caractéristiques de l'image
- **Couches de récurrence** : étude séquentielle des caractéristiques
- **Algorithme CTC** : fonction de coût pour entraîner le modèle en associant les sorties du réseau à la vérité terrain

### Problème

- Impossible de paralléliser les calculs à cause de la séquentialité des couches de récurrence



# Fully-Convolutional Networks (FCN)

## Principe

**Supprimer la partie récurrente** au profit de couches de convolution

### Avantages

- **Parallélisation** entraînement
- **Résultats à l'état de l'art**

### Inconvénients

- Architectures très **complexes**
  - Plein de mécanismes combinés complexes (attention, normalisation, ...)
  - Mise au point expérimentale non triviale

[M. Yousef et al. (2020), D. Coquenat et al. (2020)]

# Transformers [Vaswani et al. (2017)]

**Émergence** dans nombreux domaines :

- Reconnaissance de parole
- Traitement du langage
- ...

## Avantages

- **Parallélisation** de l'entraînement
- **Analyse séquentielle** des données
- Mécanisme d'**Attention**

## Inconvénients

- **Besoin** de **nombreuses données** pour l'entraînement

[Kang et al. (2020), Barrere et al. (2022), Wick et al. (2022)]

# Faire face au manque de données

- **Génération de données**

- GAN (Generative Adversarial Network) [S. Fogel et al. (2020)]
- Utilisation de polices manuscrites

# Faire face au manque de données

- Génération de données
  - GAN (Generative Adversarial Network) [S. Fogel et al. (2020)]
  - Utilisation de polices manuscrites
- **Augmentation de données** [Y. Soullard et al. (2019)]
  - Transformations (rotation, scaling, projection, ...)

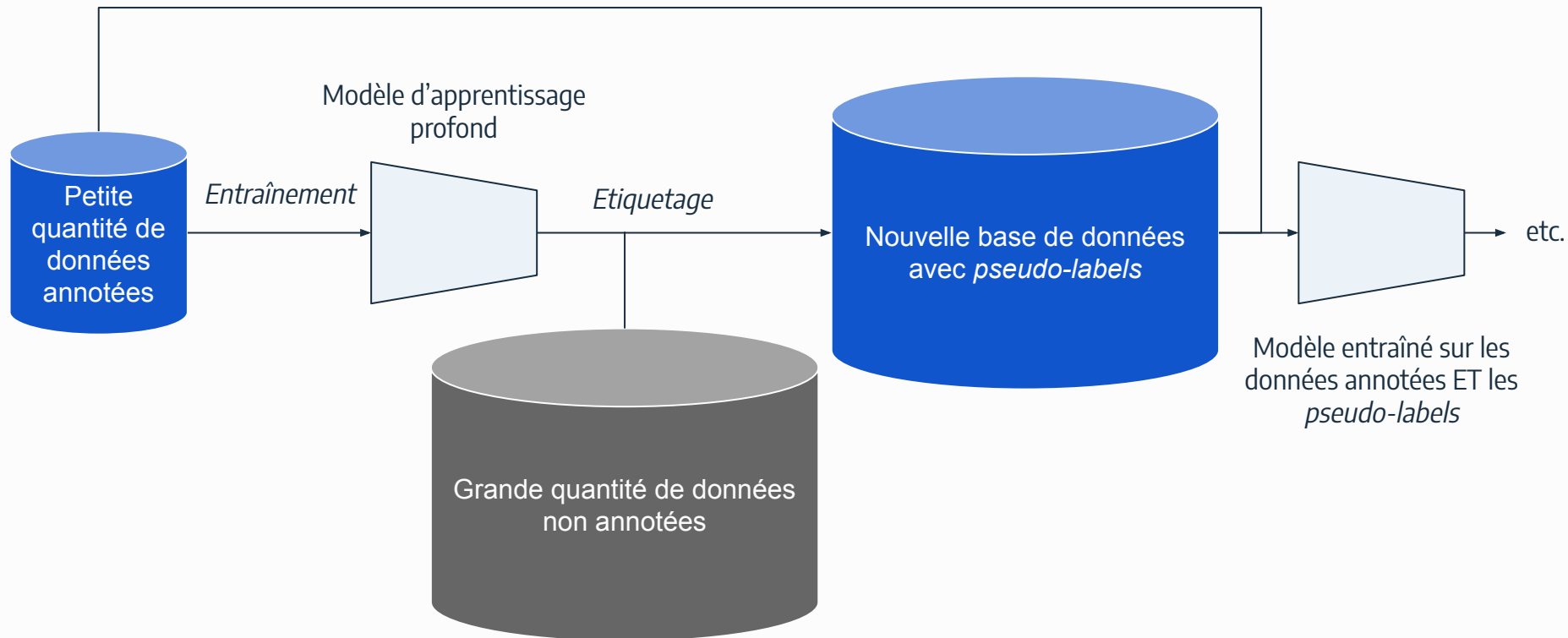
# Faire face au manque de données

- Génération de données
  - GAN (Generative Adversarial Network) [S. Fogel et al. (2020)]
  - Utilisation de polices manuscrites
- Augmentation de données [Y. Soullard et al. (2019)]
  - Transformations (rotation, scaling, projection, ...)
- **Approches semi-supervisées**
  - Faire usage des données non étiquetées

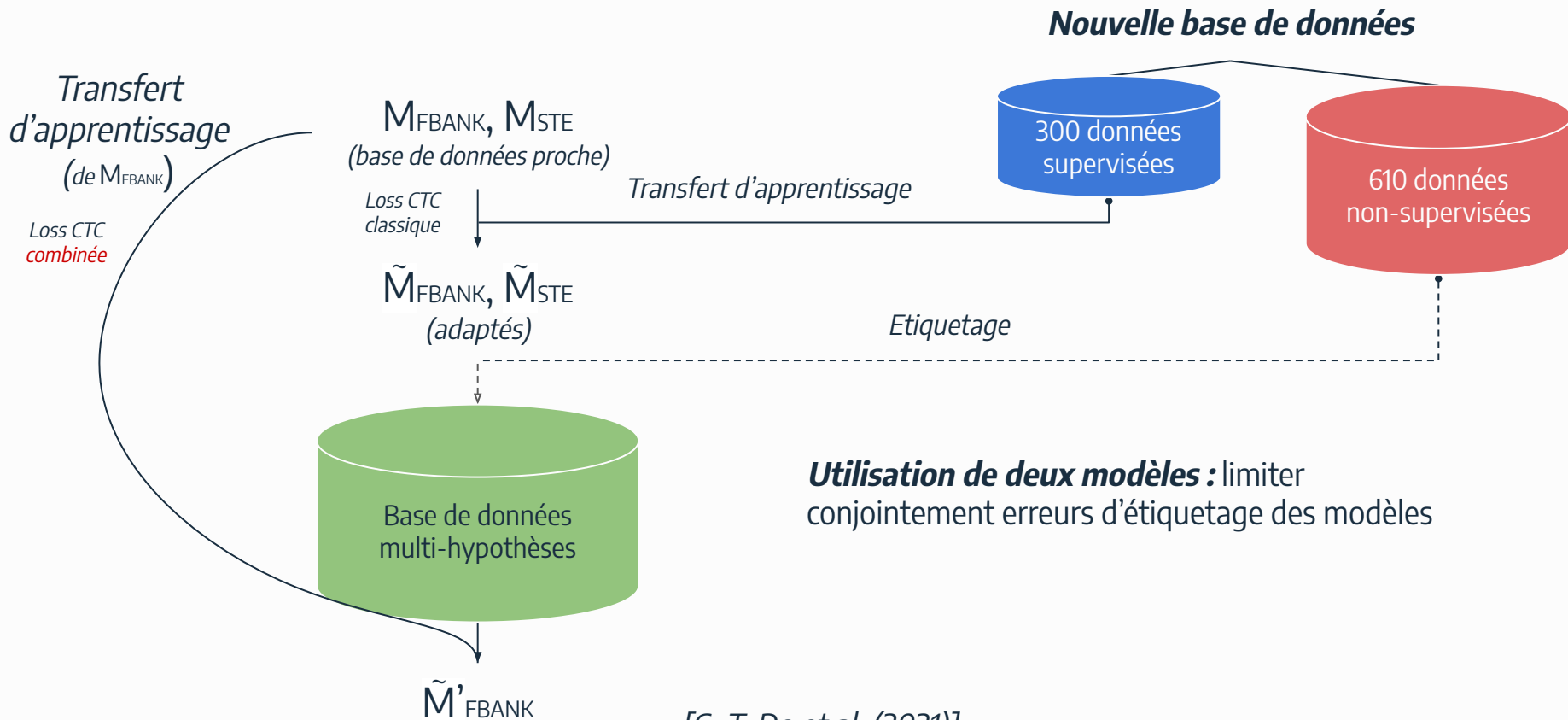
⇒ Peu de littérature en reconnaissance d'écriture manuscrite

# Approches semi-supervisées

Une façon de faire : produire des pseudo-labels sur les données non annotées



# Semi-supervisé multi-hypothèses : parole

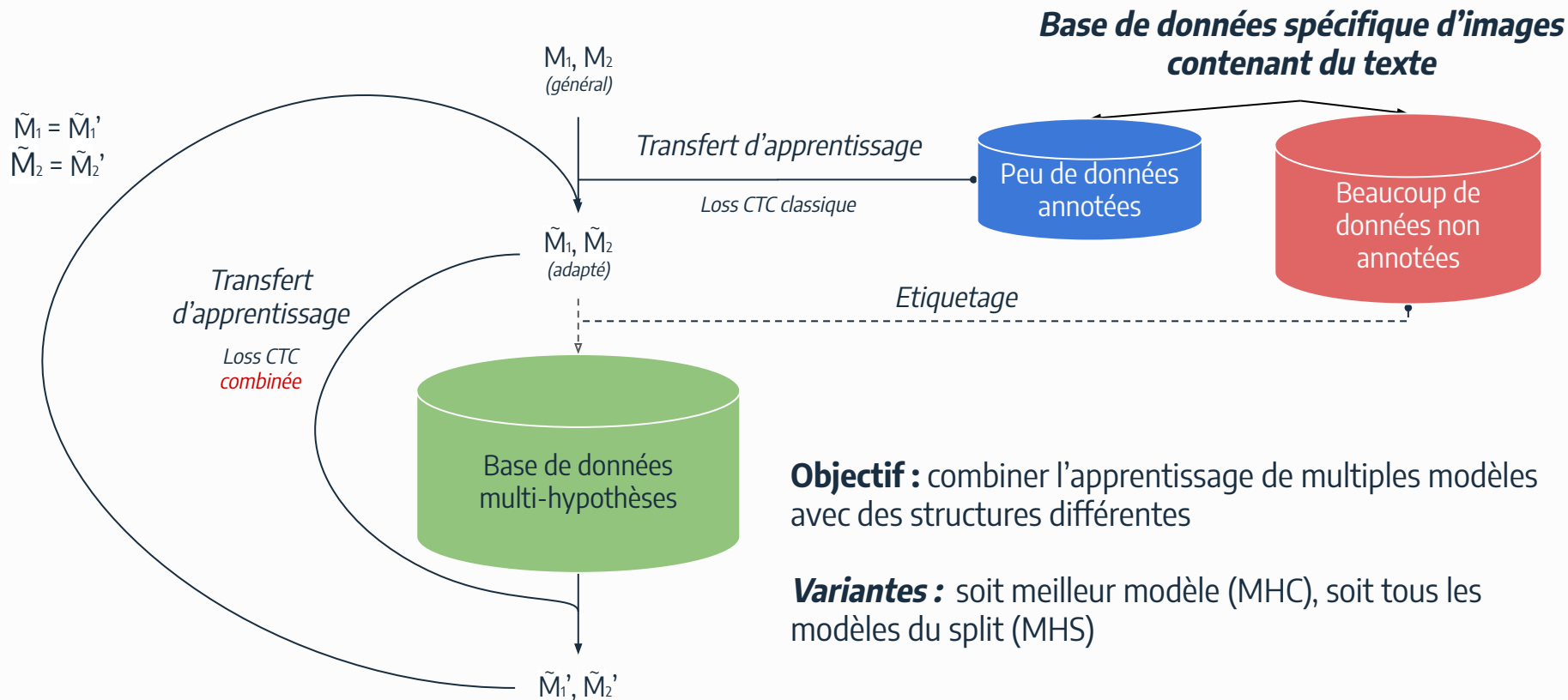


02

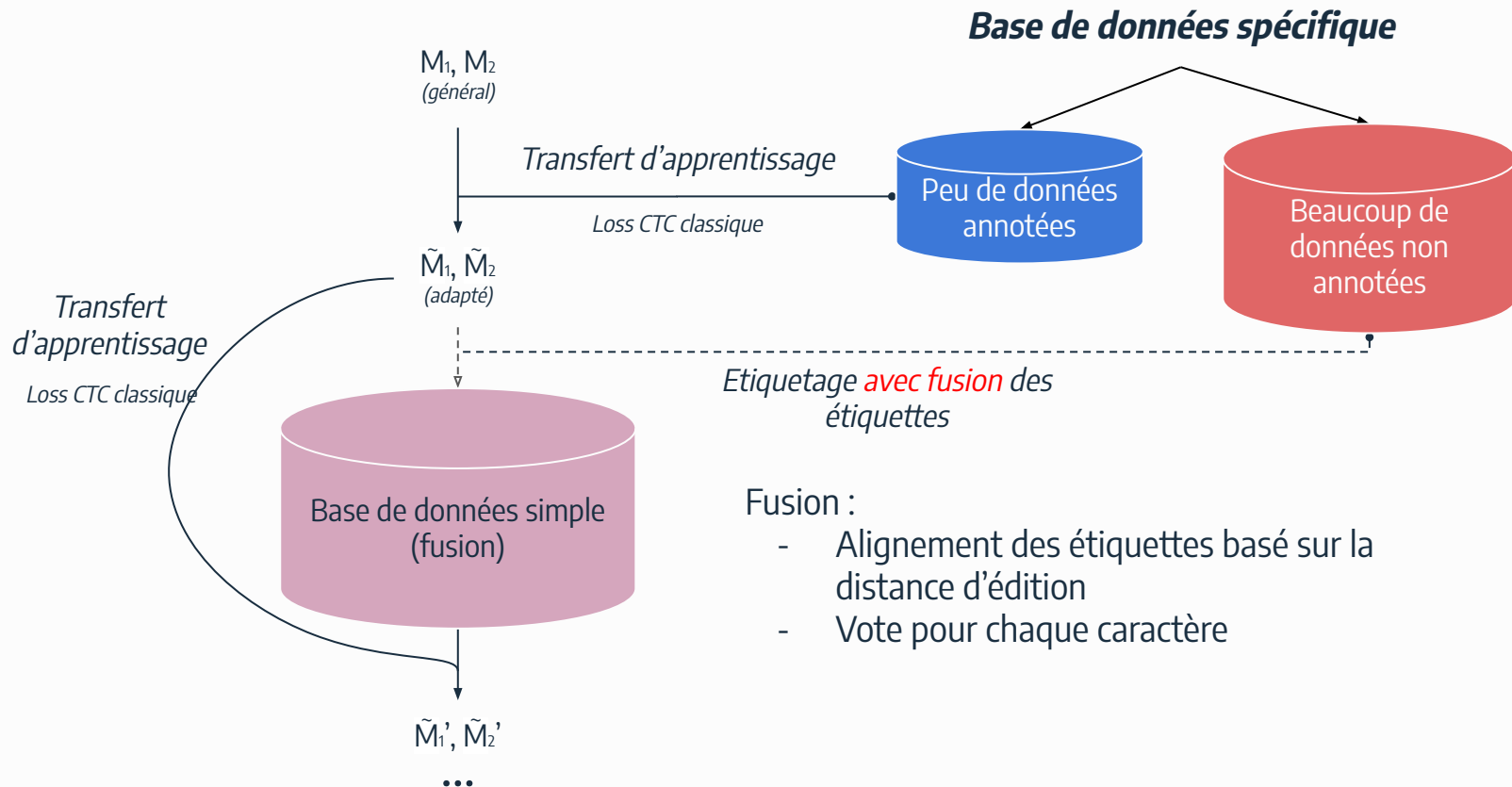
**Contribution**



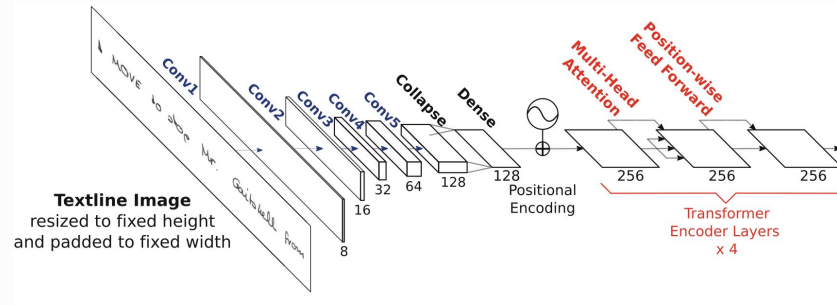
# Semi-supervisé multi-hypothèses : principe



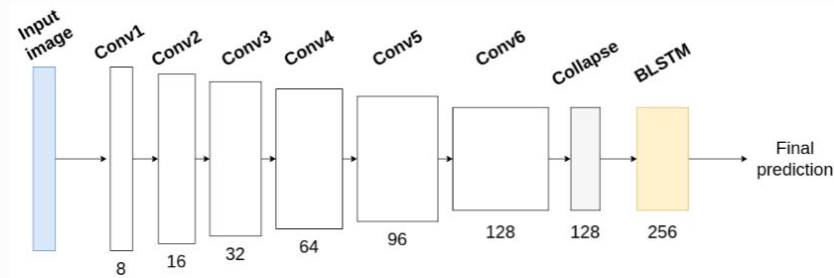
# Variante : Stratégie avec fusion



# Choix des modèles



CTNN léger (Transformer) → M1  
[K. Barrere et al. (2022)]



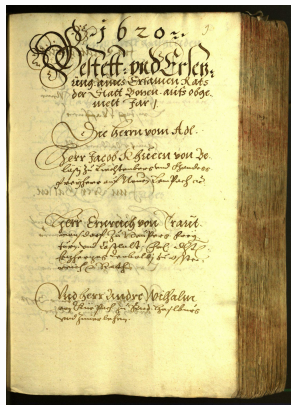
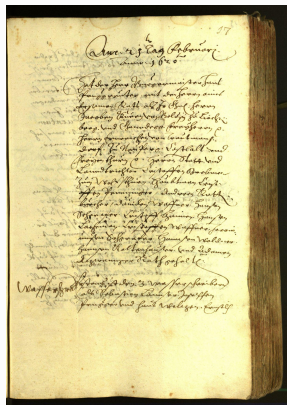
CRNN classique → M2

03

Expériences

# Base de données READ18

ICFHR 2018 : compétition “Automated Text Recognition on a READ dataset”



	general data	specific data	test data	total
documents	17	5	5	22

Différentes langues : Allemand moderne et médiéval, Anglais, Danois, Italien,...

**Character Error Rate (CER)** : taux de caractères incorrectement reconnus par rapport au texte de référence (**distance de Levenshtein**).

Plus le CER est faible, mieux c'est !

## Objectif :

- Pré-entraînement sur *general data*
- Obtenir un CER minimum sur *specific data*
- Vérifier l'impact du nombre de données avec : 1, 4 et 16 pages annotées utilisées

# Quelle stratégie adopter ? Et dans quelle configuration ?

Fusion (FUS) ? Multi-Hypothèses Choix (MHC) ?  
Multi-Hypothèses Split (MHS) ?

# Comparaison des stratégies

0 page for adaptation					
	NoA <sup>1</sup>	SA <sup>2</sup>	MHC	MHS	FUS
Konzils.	17.10	-	15.15	<b>13.97</b>	16.22
Ricordi	29.27	-	25.94	26.45	<b>24.46</b>
Schiller	<b>31.33</b>	-	41.44	100	34.57
Patzig	35.99	-	40.08	42.01	<b>35.45</b>
Schwerin	33.03	-	33.05	<b>26.25</b>	26.68
Mean $M_1$	29.35	-	29.77	41.73	<b>27.48</b>
Konzils.	19.04	-	10.54	<b>9.55</b>	11.59
Ricordi	43.96	-	23.81	24.09	<b>22.48</b>
Schiller	33.79	-	41.13	100	<b>32.94</b>
Patzig	47.97	-	36.40	39.00	<b>29.96</b>
Schwerin	53.25	-	26.73	<b>25.33</b>	26.06
Mean $M_2$	39.60	-	27.72	39.67	<b>24.61</b>

<sup>1</sup>NoA: No adaptation

<sup>2</sup>SA: Simple adaptation

- Meilleure stratégie : **FUS**
- Le **CRNN** semble d'avantages profiter de la procédure que le **CTNN**
- FUS semble capable de limiter les erreurs de chaque modèle

1 page for adaptation					
	NoA <sup>1</sup>	SA <sup>2</sup>	MHC	MHS	FUS
Konzils.	17.10	<b>10.96</b>	<b>10.89</b>	12.81	13.32
Ricordi	29.27	22.36	<b>20.97</b>	<b>21.05</b>	22.56
Schiller	31.33	<b>23.04</b>	26.74	27.89	26.61
Patzig	35.99	<b>25.67</b>	27.15	30.09	31.24
Schwerin	33.03	12.29	9.54	<b>8.75</b>	9.63
Mean $M_1$	29.35	<b>18.86</b>	19.06	20.12	20.67
Konzils.	17.10	12.38	9.06	<b>8.53</b>	9.46
Ricordi	29.27	28.57	19.04	<b>17.70</b>	19.68
Schiller	31.33	23.76	<b>20.35</b>	20.95	<b>20.37</b>
Patzig	35.99	29.08	27.15	24.90	<b>23.01</b>
Schwerin	33.03	13.26	<b>7.92</b>	9.37	9.06
Mean $M_2$	39.26	21.41	16.70	<b>16.29</b>	<b>16.32</b>

<sup>1</sup>NoA: No adaptation

<sup>2</sup>SA: Simple adaptation

- CRNN, CTNN : **MHS, SA**
- **CRNN** améliore grandement ses résultats grâce à notre méthode **MHS**
- Impact de la **fusion** semble **moins important** étant donné que les modèles sont adaptés

4 pages for adaptation					
	NoA <sup>1</sup>	SA <sup>2</sup>	MHC	MHS	FUS
Konzils.	17.10	<b>7.82</b>	<b>7.81</b>	9.54	9.53
Ricordi	29.27	<b>17.69</b>	18.50	18.85	19.66
Schiller	31.33	<b>20.10</b>	<b>20.08</b>	21.70	23.99
Patzig	35.99	<b>20.89</b>	<b>20.90</b>	24.38	21.13
Schwerin	33.03	8.83	<b>7.33</b>	<b>7.22</b>	8.80
Mean $M_1$	29.35	<b>15.07</b>	<b>14.92</b>	16.34	16.62
Konzils.	17.10	8.13	<b>7.14</b>	7.32	<b>7.03</b>
Ricordi	29.27	21.42	16.43	<b>15.87</b>	<b>16.01</b>
Schiller	31.33	20.03	16.60	<b>16.18</b>	18.03
Patzig	35.99	21.64	19.88	<b>19.10</b>	<b>19.19</b>
Schwerin	33.03	7.25	5.84	<b>5.37</b>	5.65
Mean $M_2$	39.26	15.69	13.18	<b>12.77</b>	13.18

<sup>1</sup>NoA: No adaptation

<sup>2</sup>SA: Simple adaptation

**Comment se comportent nos  
stratégies face à une approche  
semi-supervisée classique ?**



# Comparaison à du semi-supervisé classique

0 page for adaptation			
	NoA <sup>1</sup>	Semi-sup.	FUS
Konzils.	17.10	16.78	<b>16.22</b>
Ricordi	29.27	29.15	<b>24.46</b>
Schiller	31.33	<b>31.16</b>	34.57
Patzig	35.99	35.82	<b>35.45</b>
Schwerin	33.03	33.05	<b>26.68</b>
Mean $M_1$	29.35	29.19	<b>27.48</b>
Konzils.	19.04	18.90	<b>11.59</b>
Ricordi	43.96	44.03	<b>22.48</b>
Schiller	33.79	33.64	<b>32.94</b>
Patzig	47.97	47.71	<b>29.96</b>
Schwerin	53.25	53.27	<b>26.06</b>
Mean $M_2$	39.60	39.51	<b>24.61</b>

<sup>1</sup>NoA: No adaptation

- **Notre stratégie obtient de meilleurs résultats :**  
Correction commune des erreurs

1 page for adaptation			
	NoA <sup>1</sup>	Semi-sup.	MHS
Konzils.	17.10	<b>10.05</b>	12.81
Ricordi	29.27	<b>19.85</b>	21.05
Schiller	31.33	<b>21.80</b>	27.89
Patzig	35.99	<b>24.62</b>	30.09
Schwerin	33.3	10.32	<b>8.75</b>
Mean $M_1$	29.35	<b>17.33</b>	20.12
Konzils.	17.10	9.40	<b>8.53</b>
Ricordi	29.27	25.03	<b>17.70</b>
Schiller	31.33	22.12	<b>20.95</b>
Patzig	35.99	27.15	<b>24.90</b>
Schwerin	33.03	10.96	<b>9.37</b>
Mean $M_2$	39.60	18.93	<b>16.29</b>

<sup>1</sup>NoA: No adaptation

- **CRNN** profite grandement de l'entraînement combiné
- **CTNN** obtient de meilleurs résultats seul avec une approche classique.
- Les meilleurs résultats globaux sont obtenus grâce à notre stratégie MHS dans le cadre d'un processus d'adaptation

4 pages for adaptation			
	NoA <sup>1</sup>	Semi-sup.	MHS
Konzils.	17.10	<b>7.76</b>	9.54
Ricordi	29.27	<b>16.08</b>	18.85
Schiller	31.33	<b>19.33</b>	21.70
Patzig	35.99	<b>20.39</b>	24.38
Schwerin	33.03	7.35	<b>7.22</b>
Mean $M_1$	29.35	<b>14.8</b>	16.34
Konzils.	17.10	<b>6.98</b>	7.32
Ricordi	29.27	19.33	<b>15.87</b>
Schiller	31.33	21.09	<b>16.18</b>
Patzig	35.99	21.10	<b>19.10</b>
Schwerin	33.03	6.44	<b>5.37</b>
Mean $M_2$	39.26	15.11	<b>12.77</b>

<sup>1</sup>NoA: No adaptation

# Discussion sur notre approche multi-hypothèse

- Le **CRNN bénéficie** des bonnes capacités prédictives du CTNN et s'améliore le plus.
- Le **CTNN est impacté par des prédictions imparfaites** issues du CRNN.
- La stratégie de Fusion semble utile pour **réduire les erreurs de reconnaissance** lorsque le modèle se trompe souvent (cas sans adaptation).
- Dans le cas contraire, où les pseudo-labels sont plus faiblement erronées, **l'apprentissage des multiples hypothèses sans correction par fusion semble préférable.**

04

## Conclusion

# Conclusion

- Proposition d'une **nouvelle procédure d'apprentissage semi-supervisée** pour de la reconnaissance d'écriture manuscrite ancienne lorsque **peu de données annotées** sont accessibles
- Combinaison de réseaux de neurones profonds et expérimentation de différentes stratégies
- Des **résultats encourageants et proches de l'état de l'art** sur une base de données de référence malgré des modèles initiaux pas aussi performants

## Perspectives

- Expérimenter l'approche avec des architectures à l'état de l'art
- Intégrer un modèle de langues au sein de l'approche pour apporter des corrections linguistiques aux prédictions

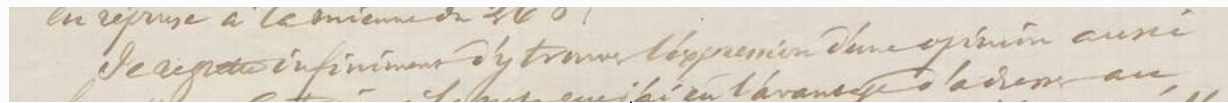
# Bibliographie

- Graves et al. (2006).  
*"Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks"*
- Strauß et al. (2018).  
*"Icfhr2018 competition on automated text recognition on a read dataset"*
- M. Yousef et al. (2020).  
*"Accurate, data-efficient, unconstrained text recognition with convolutional neural networks"*
- Vaswani et al. (2017).  
*"Attention is all you need"*
- K. Barrere et al. (2022).  
*"A light Transformer-Based Architecture for Handwritten Text Recognition"*
- L. Kang et al. (2020).  
*"Pay attention to what you read: Non-recurrent handwritten text-line recognition"*
- Y. Soullard et al. (2019).  
*"Improving text recognition using optical and language model writer adaptation"*
- E. Chammas et al. (2018).  
*"Handwriting recognition of historical documents with few labeled data"*

# Bibliographie

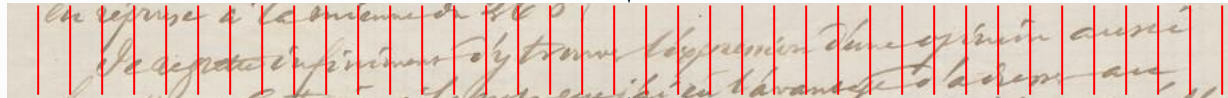
- S. Fogel et al. (2020).  
*“Scrabblegan: Semi-supervised varying length handwritten text generation”*
- C.-T. Do et al. (2021).  
*“Multiple-hypothesis ctc-based semi-supervised adaptation of end-to-end speech recognition”*
- G. Leifert et al. (2020).  
*“Two semi-supervised training approaches for automated text recognition”*
- B. Shi et al. (2016).  
*“An end-to-end trainable neural network for image-based sequence recognition and its application to scene-text recognition”*

# Modèles classiques : CRNN



Extraction des  
caractéristiques de l'image  
d'entrée

CNN



Etude des caractéristiques de  
manière séquentielle

BLSTM

Relie séquences à vérité  
terrain

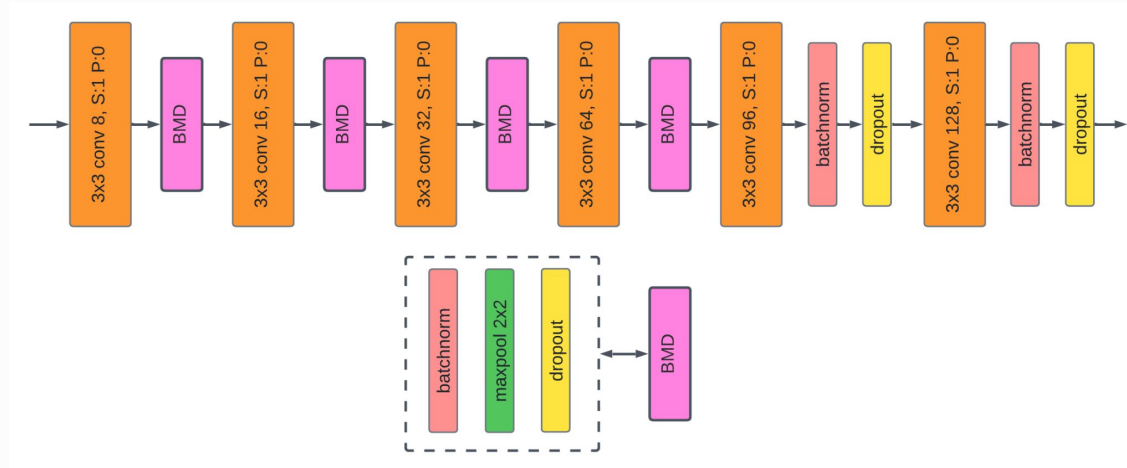
CTC

Matrice de proba.



$$\begin{matrix} & t_1 & t_2 & t_3 & \cdots & t_n \\ \begin{matrix} a \\ b \\ c \\ d \\ e \\ \vdots \\ z \\ - \end{matrix} & \begin{pmatrix} 0.1 & 0.0 & 0.0 & \cdots & 0.01 \\ 0.5 & 0.0 & 0.01 & \cdots & 0.1 \\ 0.0 & 0.0 & 0.4 & \cdots & 0.0 \\ 0.0 & 0.2 & 0.1 & \cdots & 0.4 \\ 0.01 & 0.3 & 0.0 & \cdots & 0.01 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0.1 & 0.1 & 0.2 & 0.0 & 0.02 \\ 0.02 & 0.11 & 0.12 & 0.0 & 0.1 \end{pmatrix} \end{matrix}$$

# Choix des modèles



CRNN inspiré par celui de l'Université  
de l'Ohio  
[B. Shi et al. (2016)]



# Modèles classiques : Loss CTC

$$\begin{array}{c}
\text{d} \left\{ \begin{array}{c}
\begin{array}{c}
a \\ b \\ c \\ d \\ e \\ \vdots \\ z \\ -
\end{array}
\begin{pmatrix}
\begin{array}{ccccc}
t_1 & t_2 & t_3 & \cdots & t_n \\
0.1 & 0.0 & 0.0 & \cdots & 0.01 \\
0.5 & 0.0 & 0.01 & \cdots & 0.1 \\
0.0 & 0.0 & 0.4 & \cdots & 0.0 \\
0.0 & 0.2 & 0.1 & \cdots & 0.4 \\
0.01 & 0.3 & 0.0 & \cdots & 0.01 \\
\vdots & \vdots & \vdots & \ddots & \vdots \\
0.1 & 0.1 & 0.2 & 0.0 & 0.02 \\
0.02 & 0.11 & 0.12 & 0.0 & 0.1
\end{array}
\end{pmatrix}
\end{array}
\right.
\end{array}$$

Soit une **séquence de features** de T caractères :

$$X = \{x_t \in \mathbb{R}^d | t = 1, \dots, T\}$$

Soit une **transcription** de L caractères :

$$C = \{c_l \in U | l = 1, \dots, L\}$$

La **fonction de coût CTC** est calculée en intégrant sur tous les chemins qui peuvent être associés à  $C$ ,  $B^{-1}(C)$ :

$$L_{CTC} = -\log \sum_{\alpha \in B^{-1}(C)} P_{\theta}(\alpha | X)$$

# Semi-supervisé multi-hypothèses : loss combinée

$\hat{C}_i, i = 1, \dots, N$  : les **pseudo-labels** générés par les  $N$  modèles ( $N = 2$  dans le papier)

Redéfinition de la **fonction de coût CTC** :

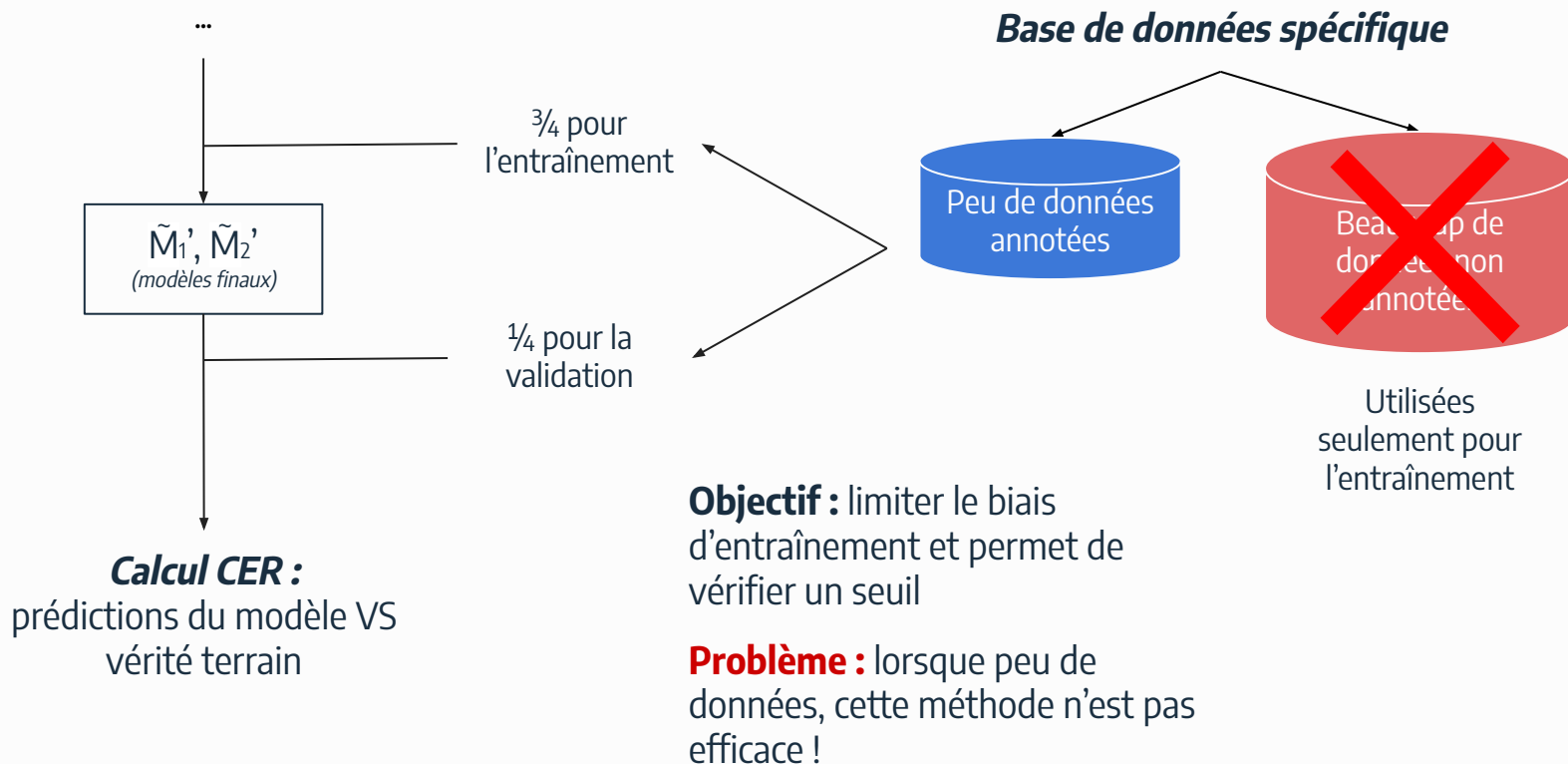
$$L_{CTC}^* = - \left( \sum_{i=1}^N \log P_{\theta}(\hat{C}_i | X) \right)$$

Avec les propriétés de la fonction *log* on peut écrire :

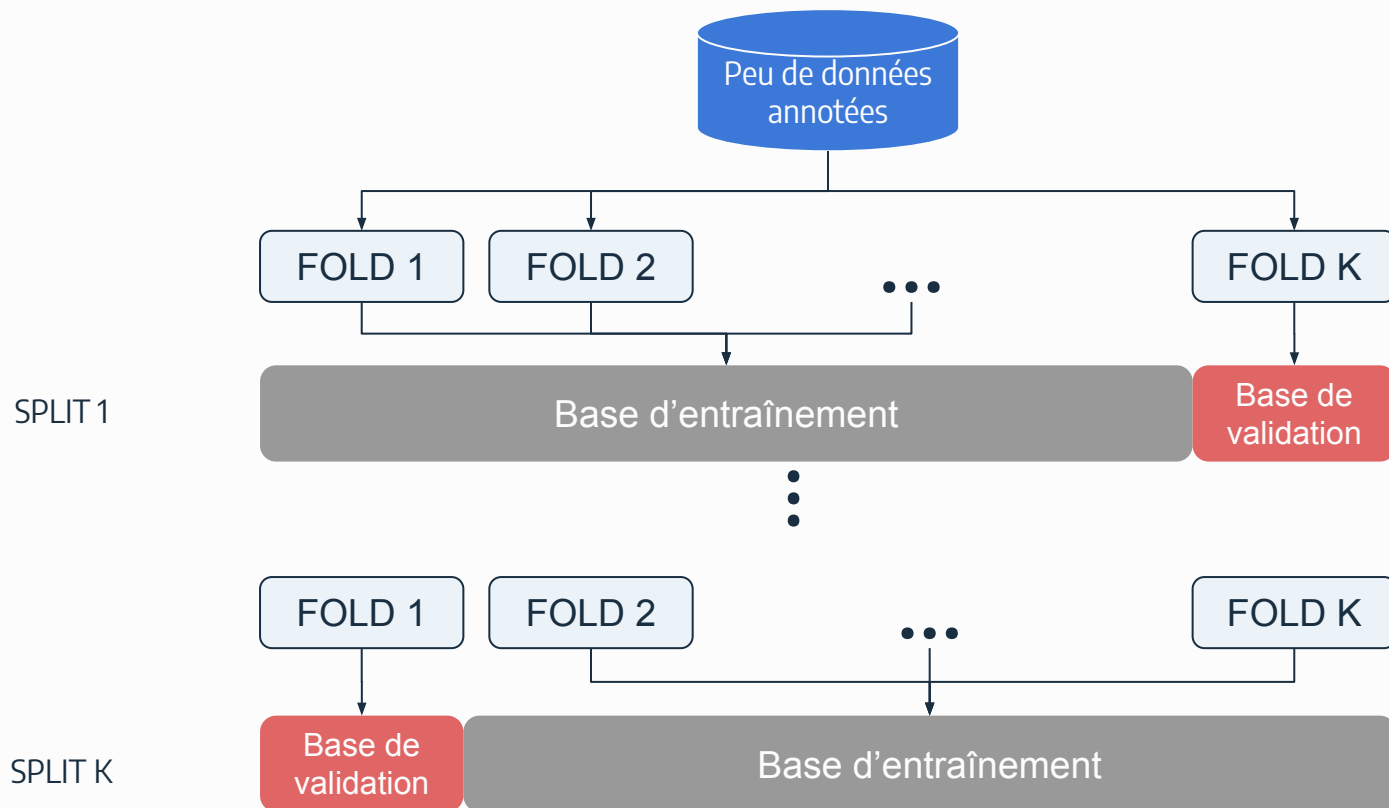
$$L_{CTC}^* = - \log \left[ \left( \sum_{a_i \in B^{-1}(\hat{C}_1)} P_{\theta}(a_i | X) \right) \left( \sum_{b_i \in B^{-1}(\hat{C}_2)} P_{\theta}(b_i | X) \right) \right]$$

[C.-T. Do et al. (2021)]

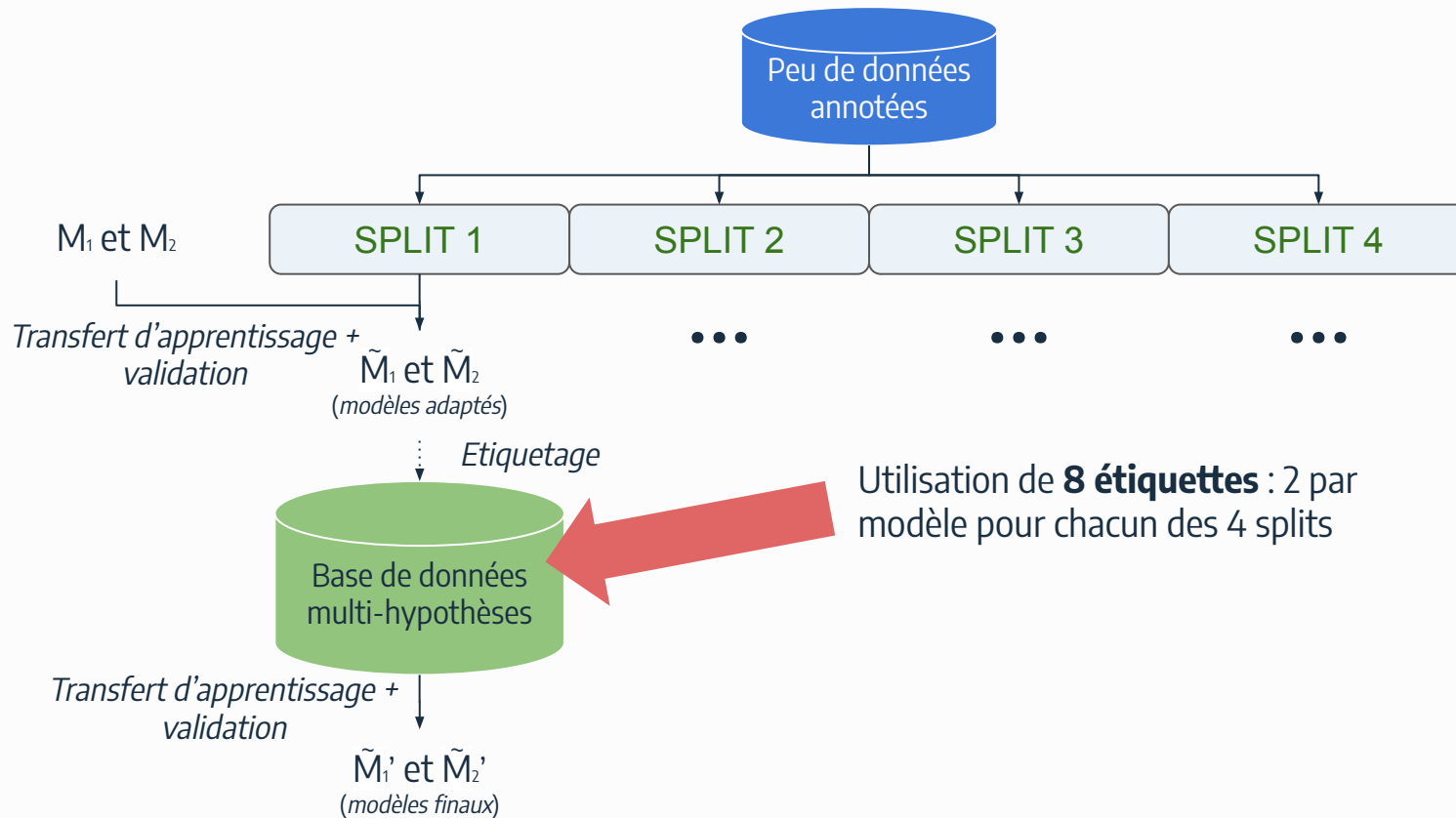
# Validation



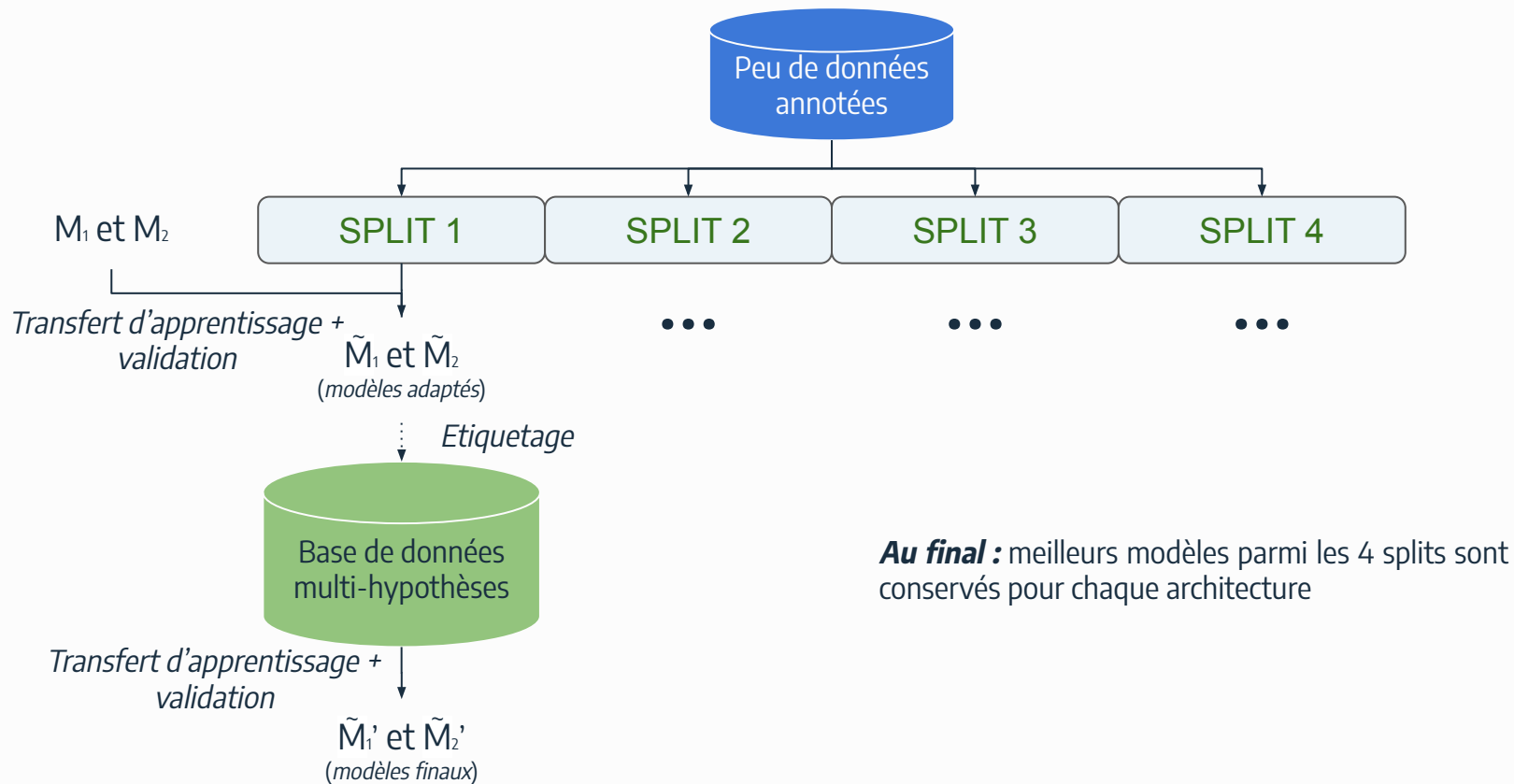
# Validation croisée (*k-fold cross validation*)



# Variante : Combinaison multi-split



# Seuil par validation croisée (4-fold)

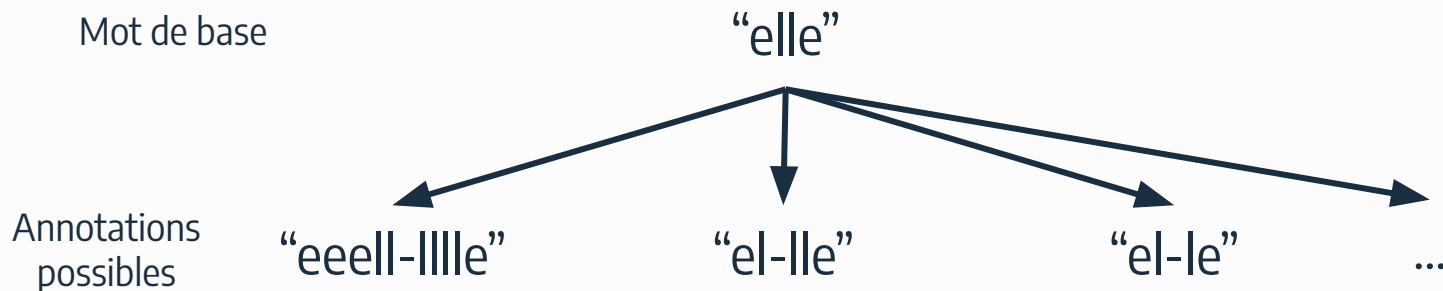


# Algorithme CTC : Decodage

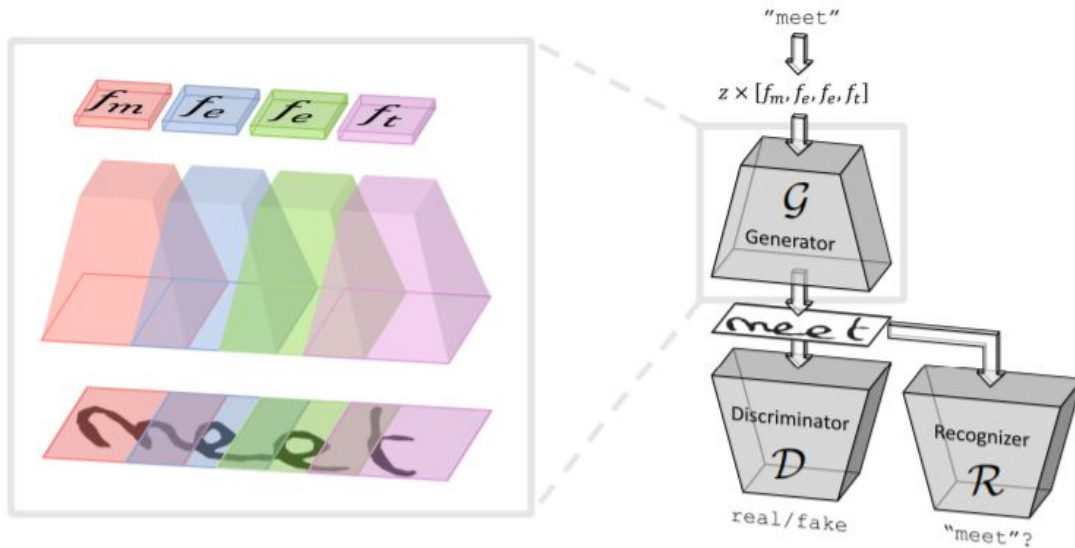
**Problème** : mots qui ont des caractères qui se répètent : *'belle', 'bille', etc.*



**Solution** : ajout du caractère "**blank**" (noté '\_') entre caractères qui se répètent



# Faire face au manque de données : génération de données



Architecture ScrabbleGan utilisée



Ensembles de mots générés par ScrabbleGan en suivant un certain style sur chaque ligne



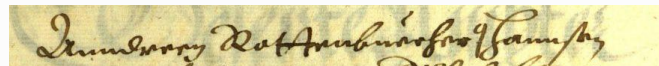
# Algorithme CTC

**Modèle** à entraîner

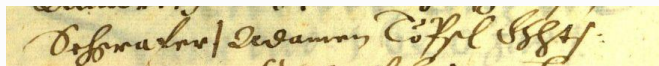


- Séquence de *features* de nombre inconnu (**non-segmentée**)

**Annotation** du document  
manuscrit



'Anndreen Rottenbuecher. Hannsen'



'Schgrafer/ Adamen TōPsl Rhts→'

- Séquence de caractères finie (**segmentée**)

**Comment aligner les deux ?**

Grâce à l'algorithme CTC !

# Algorithme CTC : Decodage

	$t_1$	$t_2$	$t_3$	$\dots$	$t_n$
$a$	0.1	0.0	0.0	$\dots$	0.01
$b$	0.5	0.0	0.01	$\dots$	0.1
$c$	0.0	0.0	0.4	$\dots$	0.0
$d$	0.0	0.2	0.1	$\dots$	0.4
$e$	0.01	0.3	0.0	$\dots$	0.01
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$z$	0.1	0.1	0.2	0.0	0.02
$-$	0.02	0.11	0.12	0.0	0.1

$\text{argmax}(t_1) = 1$   
→ caractère 'b'

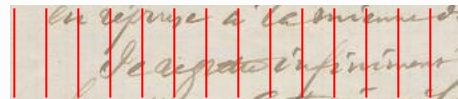
...

$\text{argmax}(t_n) = 3$   
→ caractère 'd'

⇒

Taille n

bec...d



Caractères qui se répètent car ils  
apparaissent dans plusieurs  
“tranches”

⋮

**Solution :** supprimer les  
caractères qui se répètent

⋮

**Problème :** mots qui ont  
des caractères qui se  
répètent : ‘belle’, ‘bille’, etc.

# Couches de convolution

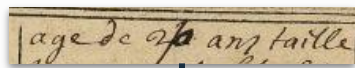
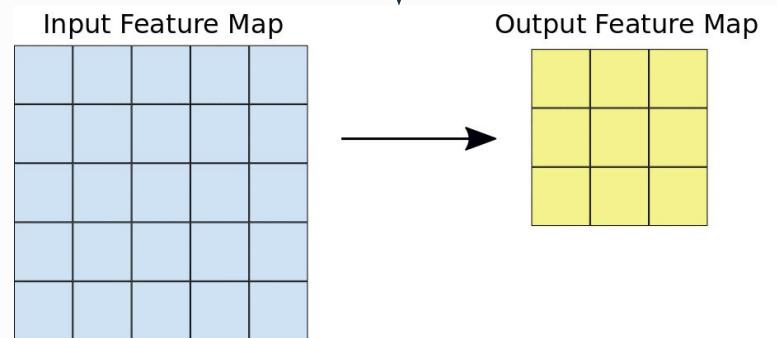
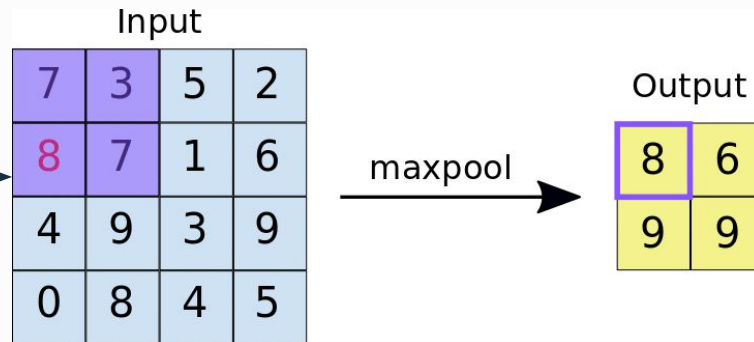


Image d'entrée



Opération de **convolution**

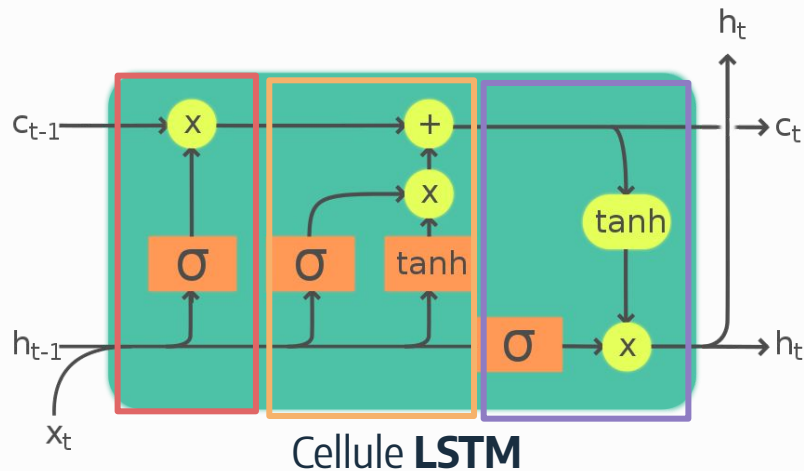
- Fenêtre glissante
- Application d'un noyau permettant d'extraire différents types de *features*



Opération de **max pooling**

- Fenêtre glissante
- Permet d'agréger localement les informations

# Couches récurrentes : LSTM et BLSTM



**BLSTM** 2 modèles LSTM apprennent la séquence en parallèle, un à l'endroit, l'autre à l'envers.

## Idées

- Séquence d'entrée  $x$
- Cell state  $c$  permet de conserver ce qui est important à long terme
- Hidden state  $h$  permet de conserver ce qui est important à court terme
- **Input gate** : ajoute informations pertinentes pour long-terme à  $c$
- **Forget gate** : détermine informations pertinentes venant du passé
- **Output gate** : récupère dans  $c$  les informations pertinentes à court terme

# Extensions possible

- Ajouter une surcouche de modèle de langue dans la procédure afin de corriger encore d'avantages les erreurs → pour se rapprocher de l'état de l'art ?
- Tenter d'autres combinaisons de modèles ? (Transformer, FCN, etc.)