



# Participation d'EDF R&D à DEFT 2022

P. Suignard, X. Huang, M. Bothua

Atelier du 27/06/2022



# Plan de la présentation

- EDF R&D
- Pourquoi participer à ce concours ?
- Les tâches et les méthodes utilisées
- Résultats obtenus
- Conclusion



- Structure

- ✓ EDF R&D au service de toutes les entités du groupe EDF
- ✓ Basée à Saclay, le nouveau centre de recherche
- ✓ Environ 2000 personnes



## ■ Travaux sur le texte

- ✓ En appui aux différents métiers : EDF Commerce, Hydraulique, Eolien, Nucléaire, Enedis, RH, IT...
- ✓ En permanence : ~8 personnes, 1 doctorant, 2 stagiaires
- ✓ Thèmes : classification, clustering, orthographe, annotations, web sémantique, résumé, anonymisation de données, détection de nouveauté, plongements mots/documents...
- ✓ Sujets : mails, réclamations, comptes-rendus d'interventions, documents techniques, manuscrits anciens, conversations téléphoniques, réseaux sociaux, chatbots...
- ✓ Type de prestations : veille, développement, conseil, méthode, étude.

# Pourquoi participer à ce concours ?

## DEFT 2022

Défi Fouille de Textes@TALN 2022

Correction automatique de copies d'étudiants

- Participation à la tâche de **base** et à la tâche **continue**
- La problématique abordée : calcul de **similarité entre paires de phrases**
- Permet de se comparer/partager/discuter avec les autres équipes
- Emulation interne
- Reconnaissance interne
- Les résultats contribuent directement à EDF Commerce et à d'autres entités du groupe EDF

# Tâche de base : Evaluation automatique de copies d'après une référence existante

- Objectif : noter des copies
- Entraînement (données de 2021)
  - ✓ On dispose de :
    - Question – q
    - Réponse attendue par l'enseignant – ra
    - Réponse de l'élève – r
    - Note obtenue
- Test
  - ✓ On dispose de q, ra et r
  - ✓ Il faut trouver la note

Qu'est-ce que le World Wide Web ?

- Système hypertexte fonctionnant sur internet
- Une des applications d'internet, comme courrier électronique, messagerie instantanée...

Ce sont les pages web accessible par tout navigateur => 0,5

Le réseau mondial, Internet => 0

C'est le systeme hypertexte qui sert à consulter des documents et des pages hébergés sur le réseau internet=> 1

# Tâche de base : Evaluation automatique de copies d'après une référence existante

## ■ Analyse des erreurs de notre système de 2021

### ■ 2004

- « Mettre » voix passive/active

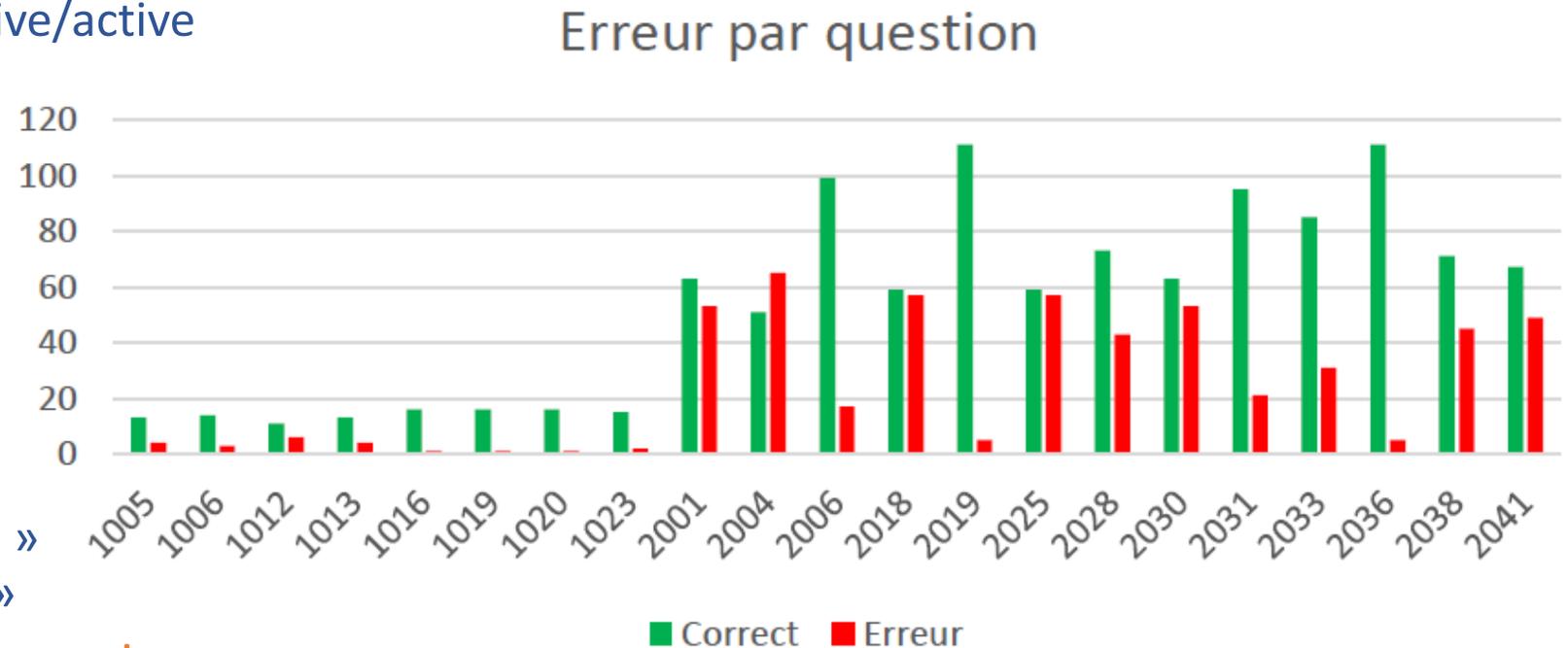
### ■ 2018

- Même réponse  
19 fois fausse et  
40 fois exacte

### ■ 2041

- Ra=« typage de contenu » alors que r=« type de données »  
« type des attributs »

- => légères modifications de ra



# Tâche de base : Evaluation automatique de copies d'après une référence existante

- Les étapes :
  - ✓ **Suppression des commentaires de l'enseignant (peuvent *perturber* le classifieur)**
    - Si ceci 1, si cela 0,5, si autre 0 => Si ceci 1
    - « Réponse attendue », toutes les réponses sont attendues !
    - « 1 si petite faute sur le texte », comment définir une « petite faute » ?
    - « 1 point pour la définition, 1 point pour l'exemple », qu'est qu'une définition, qu'est-ce qu'un exemple ?
  - ✓ **Modification**
    - changement de « i = mise en forme physique</p><p>em= mise en forme logique » par « i = mise, met, mettre en forme physique</p><p>em= mise, met, mettre en forme logique »
    - changement de « Les liens doivent décrire leur destination » par « Les liens doivent décrire leur destination, alt, alternatif »
- Méthode : Extraction de *features* (similarités croisées) + entraînement d'un classifieur

# Tâche de base : Evaluation automatique de copies d'après une référence existante

## ■ Les étapes :

### ✓ Prétraitements

- Normalisation des balises "&lt;" , "&gt;" en « < » et « > » ;
- Les balises <p>, </p>, <br> en début et fin de texte ont été supprimées ;
- Remplacement des caractères "&nbsp;" par un blanc ;
- Suppression des caractères \n et \t ;
- Insertion d'un caractère blanc avant « < » et après « > » pour faciliter la tokenisation en mots pour les calculs de similarité ;
- Utilisation du caractère blanc pour la séparation des phrases en tokens ;
- Passage en minuscule.

# Tâche de base : Evaluation automatique de copies d'après une référence existante

- Les étapes (suites):

- ✓ Extraction de 42 *features* (basées sur Soft Cardinalité, Jaro Wikler, Damereau-levenshtein et Monge Elkan)

- ✓ Soft Cardinalité (utilisé pour SemEval 2013) ou similarité floue

$$|S|' = \sum_{i=1}^n \frac{1}{\sum_{j=1}^n \text{sim}(s_i, s_j)} \quad \text{sim}(t_1, t_2) = \frac{2 * |t_1^{[2:3]} \cap t_2^{[2:3]}|}{|t_1^{[2:3]}| + |t_2^{[2:3]}|}$$

- ✓ Monge Elkan

$$\text{sim}_{\text{MongeElkan}}(S, T, \text{sim}_{\text{mot}}) = \frac{1}{|S|} * \sum_{i=1}^{|S|} \max_{j=1 \text{ à } |T|} \text{sim}_{\text{mot}}(s_i, s_j)$$

$$\text{sim}_{\text{ME}}(S, T, \text{sim}_{\text{mot}}) = \sqrt{\text{sim}_{\text{MongeElkan}}(S, T, \text{sim}_{\text{mot}}) * \text{sim}_{\text{MongeElkan}}(T, S, \text{sim}_{\text{mot}})}$$

# Tâche de base : Evaluation automatique de copies d'après une référence existante

## ■ Les étapes (suites):

### ✓ Les 42 features

- q = question
- a ou r = réponse
- ra = réponse attendue

### ✓ Entraînement d'un classifieur (Random Forest)

### ✓ Prédiction de 3 classes :

- 0 si note < 0,25
- 1 si note > 0,75
- 0,5 sinon

### ✓ +1 feature

- Cos (a, q+ra)
- sur bigrammes

$ a '$	$\frac{ a \cap ra '}{\min( a ',  ra ')}$	$\text{cosinus}_{bi}(a, ra)$
$ q '$	$\frac{ a \cap ra '}{\max( a ',  ra ')}$	$\text{cosinus}_{bi}(q, ra)$
$ ra '$	$\frac{ a \cap ra ' * ( a ' +  ra ')}{2 *  a ' *  ra '}$	$\text{cosinus}_{tri}(a, q)^1$
$ a \cup q '$	$ a \cup ra ' -  a \cap ra '$	$\text{cosinus}_{tri}(a, ra)$
$ a \cup ra '$	$\text{cosinus}_{mot}(a, q)^2$	$\text{cosinus}_{tri}(q, ra)$
$ q \cup ra '$	$\text{cosinus}_{mot}(a, ra)$	$\text{cosinus}(a, q \cup ra)$
$ a \cap ra '$	$\text{cosinus}_{mot}(q, ra)$	$\text{cosinus}_{bi}(a, q \cup ra)$
$ a \setminus ra '$	$\text{similarité}_{\text{Damereau\_Levenshtein}}(a, q)^3$	$\text{cosinus}_{tri}(a, q \cup ra)$
$ ra \setminus a '$	$\text{similarité}_{\text{Damereau\_Levenshtein}}(a, ra)$	$\text{sim}_{ME}(a, q, \text{sim}_{\text{DamereauLevenshtein}})^4$
$\frac{ a \cap ra '}{ a '}$	$\text{similarité}_{\text{Damereau\_Levenshtein}}(q, ra)$	$\text{sim}_{ME}(a, ra, \text{sim}_{\text{DamereauLevenshtein}})$
$\frac{ a \cap ra '}{ ra '}$	$\text{similarité}_{\text{JaroWinkler}}(a, q)^3$	$\text{sim}_{ME}(q, ra, \text{sim}_{\text{DamereauLevenshtein}})$
$\frac{ a \cap ra '}{ a \cup ra '}$	$\text{similarité}_{\text{JaroWinkler}}(a, ra)$	$\text{sim}_{ME}(a, q, \text{sim}_{\text{stricte}})^5$
$\frac{2 *  a \cap ra '}{ a ' +  ra '}$	$\text{similarité}_{\text{JaroWinkler}}(q, ra)$	$\text{sim}_{ME}(a, ra, \text{sim}_{\text{stricte}})$
$\frac{ a \cap ra '}{\sqrt{ a ' *  ra '}}$	$\text{cosinus}_{bi}(a, q)^6$	$\text{sim}_{ME}(q, ra, \text{sim}_{\text{stricte}})$

# Tâche de base : Evaluation automatique de copies d'après une référence existante

- Run 1 et 2 :
  - ✓ 3 classes à prédire :
    - 0 si la note était inférieure à 0,25
    - 0,5 si la note était comprise entre 0,25 et 0,75
    - 1 si la note était supérieure à 0,75.
  - ✓ Random Forest avec 100 et 200 arbres
- Run 3 :
  - ✓ classifieur entraîné à prédire la « vraie » note comprise entre 0 et 1
  - ✓ la valeur prédite a ensuite été arrondie au dixième le plus proche
  - ✓ 100 arbres Random Forest

# Tâche de base : Evaluation automatique de copies d'après une référence existante

- Résultats

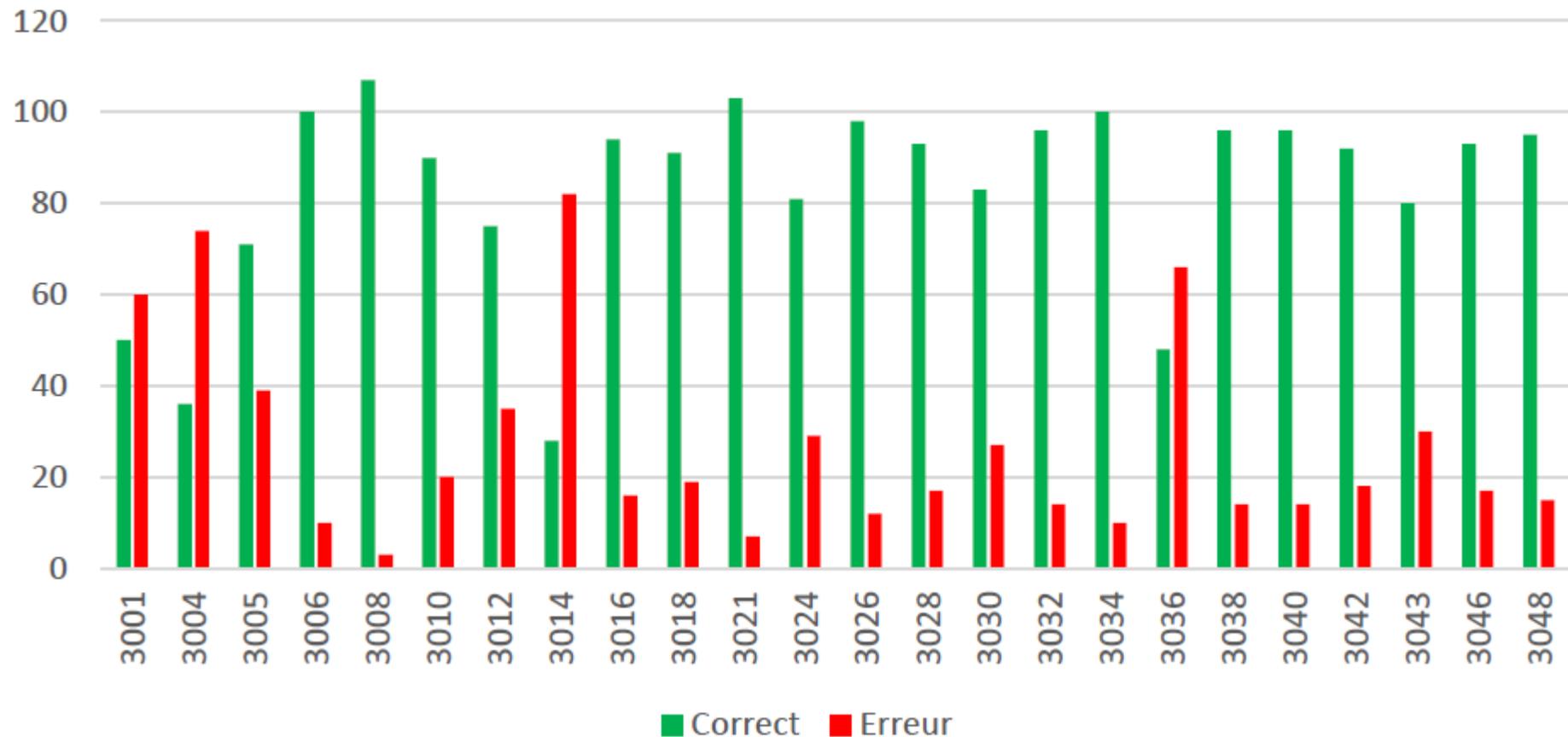
Run	Evaluation
Run 1 :	0,752
Run 2 :	0,756
Run 3 :	0,323
<i>Maximum</i>	0,756
<i>Médiane</i>	0,524
<i>Moyenne</i>	0,542
<i>Minimum</i>	0,323



# Tâche de base : Evaluation automatique de copies d'après une référence existante

- Analyse des erreurs

Erreur par question



## I. Présentation

- **Objectif**

Concevoir un système de prédiction en intégrant des interactions entre système et serveur (annotateur et oracle)

- **Données à disposition**

Les fichiers « trainT2-Q.tab » et « trainT2-R.tab » à disposition.

- **Restructuration des données:**

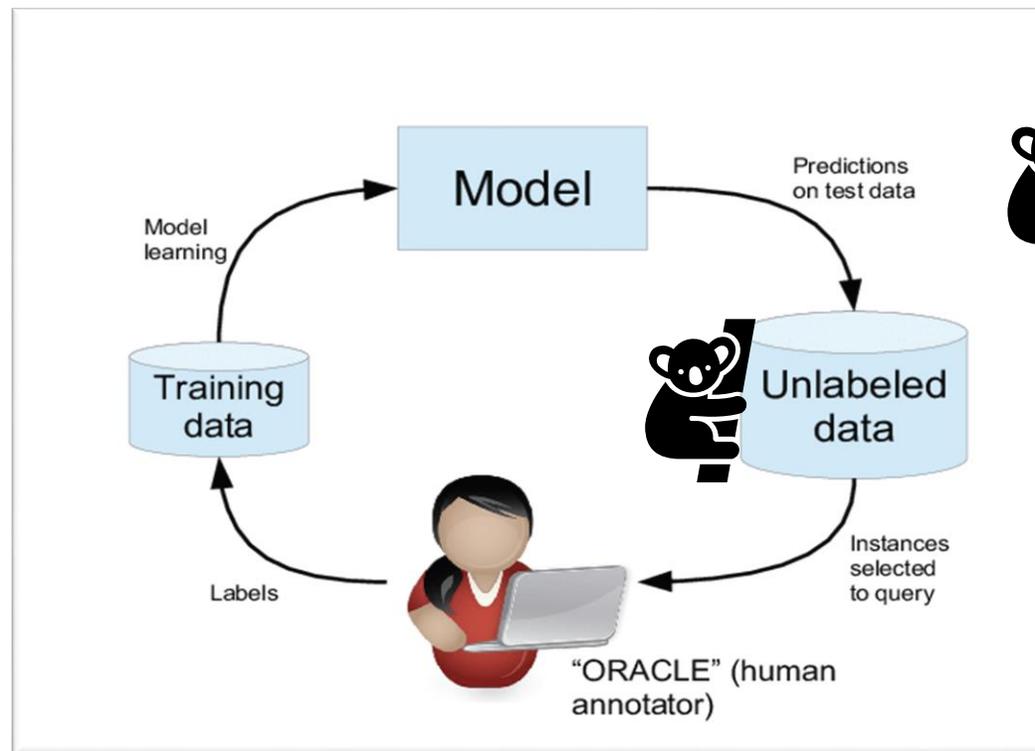
```
{'1001': [['student101', 'Ce sont les pages web accessible par tout navigateur.'],  
['student108', 'Un réseau mondial'],...]}
```

- **Etales de traitement**

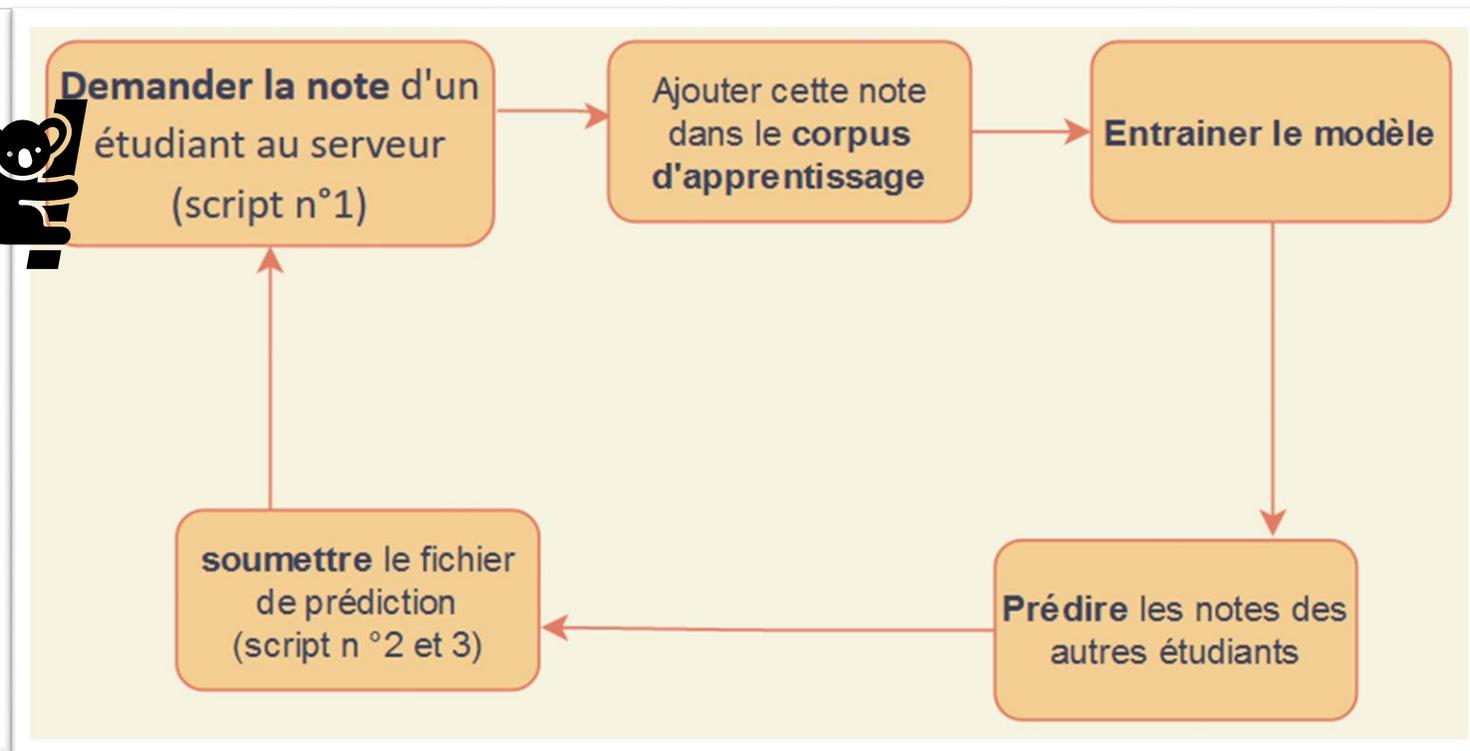
On fait référence aux principes de l'*active learning*.

# Tâche continue

## 1.3. Etapes de traitement

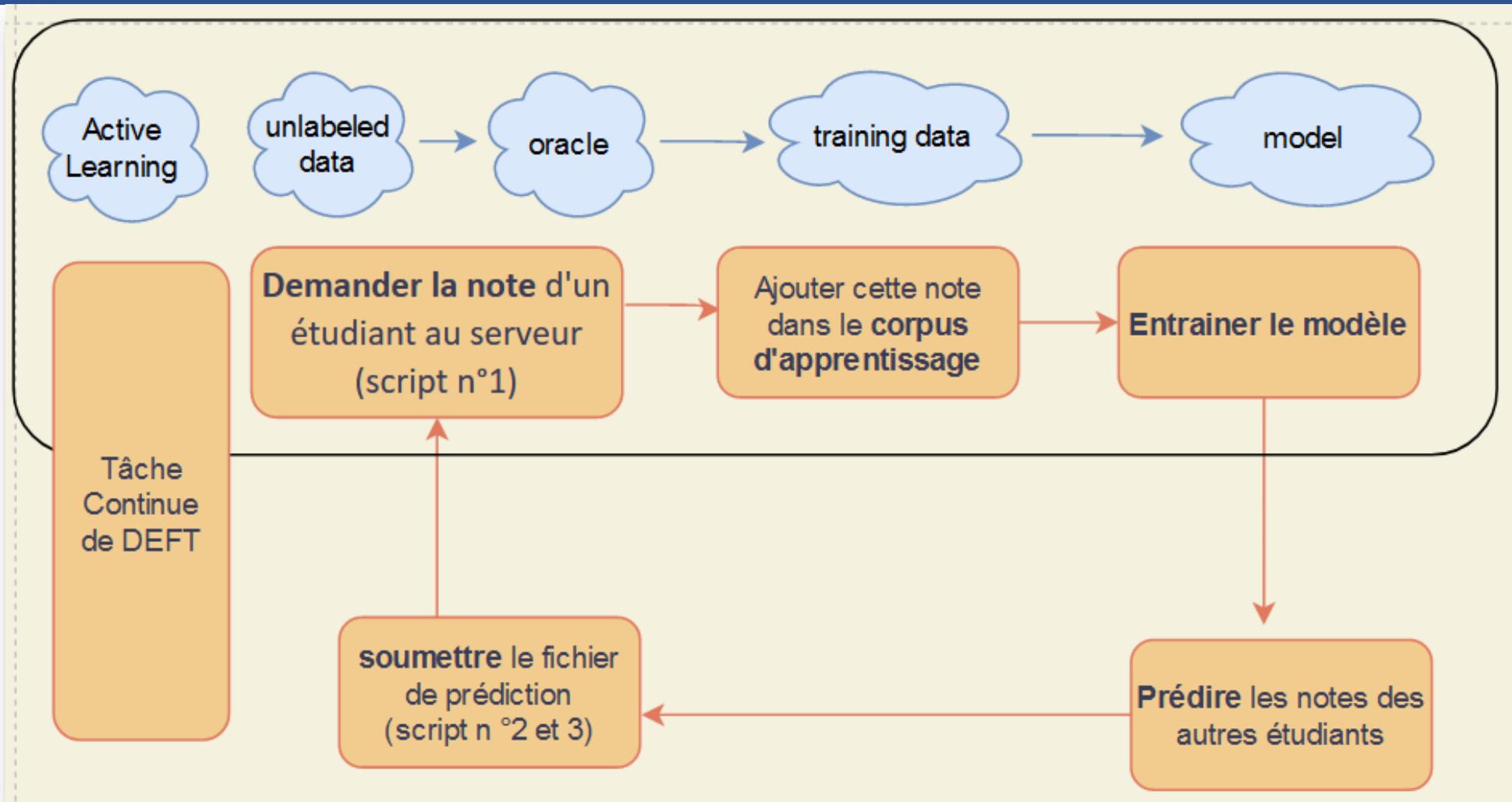


Active Learning



Tâche Continue de DEFT

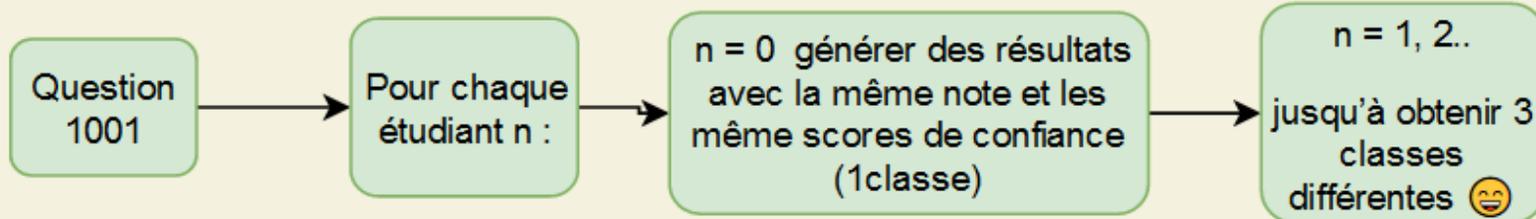
# Tâche continue



## II. Méthode

- **Méthode d'apprentissage** : 50% données pour apprentissage, 50% pour test
- **Choix de classifieur** : RandomForest
- **Feature** : Concaténation d'une réponse avec la question correspondante. Vectorisation avec SentenceTransformer.
- **Stratégie d'interrogation du serveur** : Echantillonnage d'Incertitude  
(Demander l'étiquette des données pour lesquelles le modèle actuel a le moins de certitude)

# Tâche continue



Entrée [307]:

1 train\_pool

Out [307]:

	num_ques	num_stud	grade	
0	1001	student101	0.5	Ce sont les pag
1	1001	student108	0	Un réseau mondia
2	1001	student116	0	C est le réseau
3	1001	student121	0	Le réseau mo
4	1001	student122	0	C est un reseau p
5	1001	student13	0	Le world Wide
6	1001	student3	1	C est le systèr

2002\_student68.csv

	QNUM	e_dem	etudiant	note	Confiance	sim_cos
1	2002	student68	student10	0	0.51	[tensor(0.8095)]
2	2002	student68	student100	1	0.57	[tensor(0.8723)]
3	2002	student68	student101	1	0.53	[tensor(0.8937)]
4	2002	student68	student102	0	0.6	[tensor(0.8364)]
5	2002	student68	student103	1	0.49	[tensor(0.7278)]
6	2002	student68	student104	1	0.53	[tensor(0.8725)]
7	2002	student68	student105	1	1	[tensor(0.7721)]
8	2002	student68	student106	0	0.57	
9	2002	student68	student107	1	0.59	
10	2002	student68	student107	1	0.59	

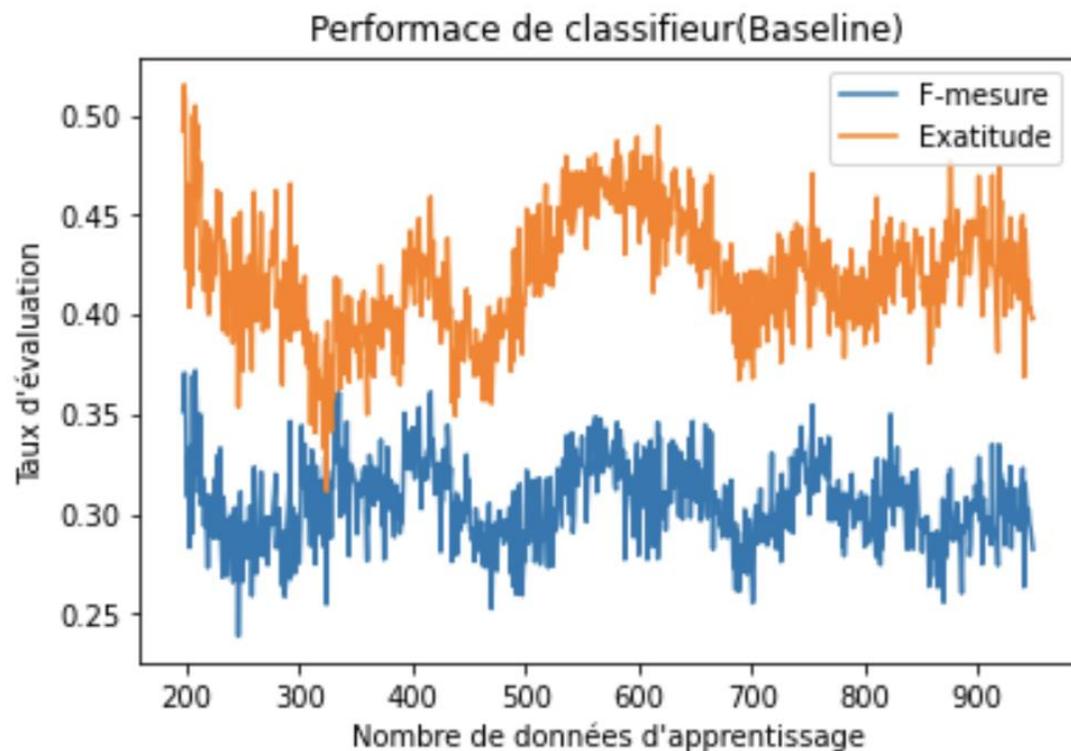
```
1001 student101 student58 0.5 1
1001 student101 student62 0.5 1
1001 student101 student7 0.5 1
1001 student101 student74 0.5 1
1001 student101 student95 0.5 1
1001 student101 student96 0.5 1
1001 student101 student98 0.5 1
```

Demander au serveur la vraie note et l'ajouter dans le corpus Train

## IV. Résultats (1)

## La performance du modèle pour la phase d'apprentissage

Classifieur : Randomforest  
Feature : Embedding: Réponse + Question



	Average	Max	Min
Accuracy	<b>0.42</b>	0.51	0.31
F-mesure	0.31	0.37	0.24

## IV. Résultats (1)

La méthode de l'apprentissage active permet de **réduire** le nombre de données d'entraînement et de **sélectionner** des données les plus représentatives (utiles).

	Nb de données	Remarques
Corpus d'apprentissage	949	24.8% du corpus complet d'apprentissage
Corpus complet d'entraînement (ensemble des fichiers train fournis)	3 820	<ul style="list-style-type: none"><li>• Question 1001 – 1028 (20 questions, 17 étudiants): 340</li><li>• Question 2002 – 2046 (30 questions, 116 étudiants): 3480</li></ul>

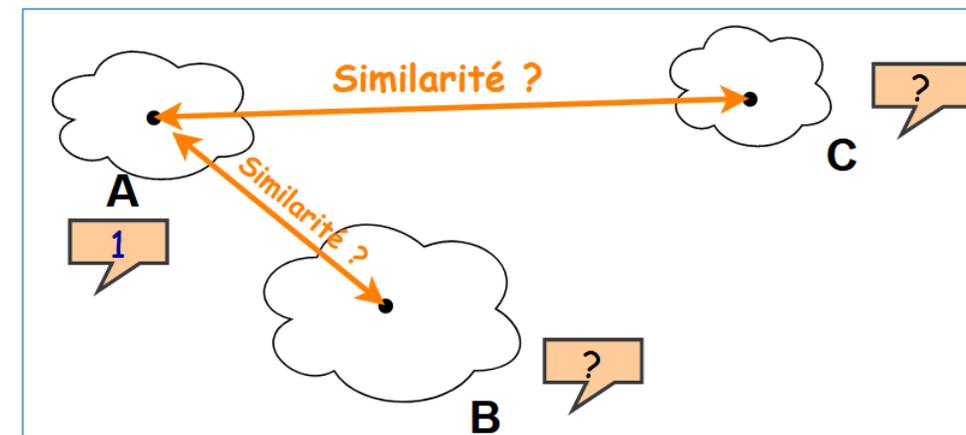
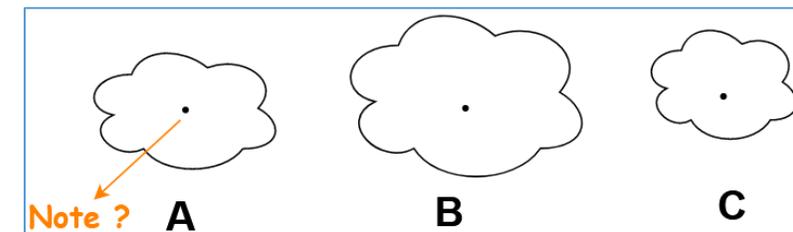
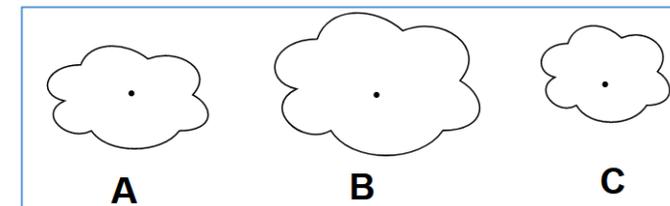
# Tâche continue

## V. Perspectives

- Réduire le nombre de fois de requête au serveur.
- Améliorer le niveau d'exactitude :
  - 1) On peut d'abord obtenir **trois clusters**.
  - 2) **Trouver** la centroïde d'un cluster et **demander** la note de ce centroïde.
  - 3) **Calculer** la similarité entre les centroïdes :

Les deux centroïdes les plus lointains auront les étiquettes 0 et 1.

- 4) Une fois que **les étiquettes des trois clusters** sont déterminées, toutes les réponses dans les clusters seront notées.



# Conclusion

- Participer à la campagne DEFT 2022, nous a permis de tester plusieurs méthodes de **calcul de similarité** et de **l'Active Learning**.
- Les résultats obtenus sont très satisfaisants, classés **premiers** sur les 2 tâches
- Les méthodes que nous avons mises en œuvre sont facilement **transposables** à d'autres tâches et peuvent intéresser plusieurs entités du groupe EDF