

Introduction aux humanités numériques – L1HN001

Nom et prénom d'étudiant-e :

Numéro d'étudiant-e :

Date d'examen : 7 décembre 2023, 08h-09h

Question 1

(4 points)

Dans la phrase suivante, soulignez toutes les entités nommées :

Emmanuel Macron, né le 21 décembre 1977 à Amiens, est président de la France

* **Remarque** : *la France* est une entité nommée de lieu qui est imbriquée à l'intérieur de l'entité nommée désignant le titre (*président de la France*).

Question 2

(1 point)

Parmi les options ci-dessous, choisissez celle(s) qui correspond(ent) à (aux) expression(s) régulière(s) valide(s) :

- ☒ ^jour
- ☐ [[^jour
- ☐ [a-zA-Z0-9_])
- ☒ [a-zA-Z0-9_]

Question 3

(1 point)

Comment s'appelle l'étape initiale d'une chaîne de traitement d'OCR ?

- ☐ exportation
- ☐ prédiction
- ☐ transcription
- ☒ segmentation

Question 4

(1 point)

Remplissez le blanc avec le mot approprié :

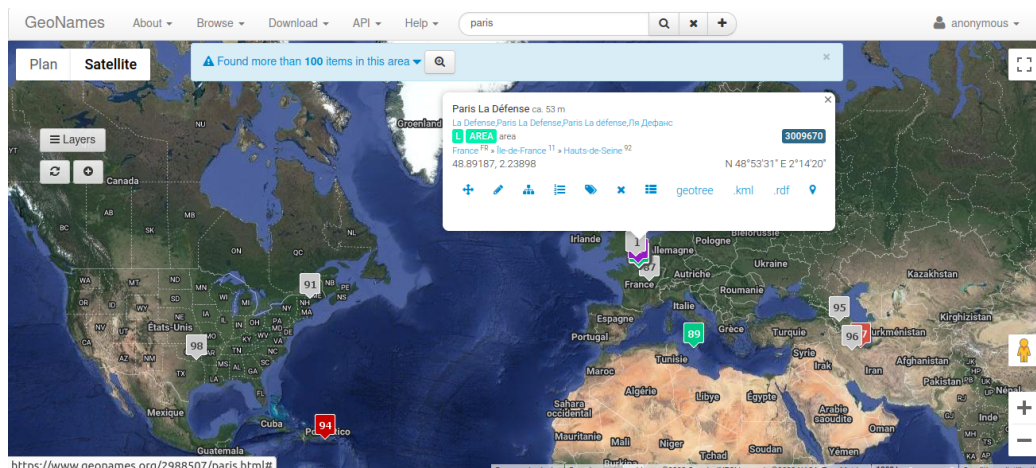
Transkribus est une plateforme de la reconnaissance / transcription automatique des textes.

Question 5

(2 points)

L'illustration suivante représente un exemple de (plusieurs réponses possibles) :

- ☐ l'annotation des entités nommées suivant le format BIO
- ☒ la cartographie des entités nommées
- ☐ la campagne d'évaluation de la reconnaissance d'entités nommées
- ☒ la désambiguïsation des entités nommées



Question 6

(1 point)

Quel est l'intérêt d'utilisation des expressions régulières ?

Les expressions régulières, aussi dénommées *regex*, fournissent un moyen concis et flexible pour la correspondance de chaînes de caractères dans un texte, telles que des caractères particuliers, mots ou motifs (patrons) de caractères.

Question 7

(2 points)

Répondez par VRAI (V) ou FAUX (F) aux phrases suivantes :

V F

- (a) L'approche à base de règles pour effectuer certaines tâches de TAL nécessite un entraînement d'un modèle. ☐ ☒
- (b) Généralement, les logiciels d'OCR / HTR permettent d'exporter des transcriptions dans différents formats. ☒ ☐

Question 8

(2 points)

Étant donné la phrase *Je suis en train de réviser pour mon 2^{ème} partiel*, formulez les regex qui permettent de capturer les mots suivants :

1. 2^{ème} : \dème ou [0-9]ème
2. Je : [A-Z]e

Question 9

(3 points)

1. Chargez le corpus *Le tour du monde en quatre-vingts jours* (utilisé lors de la séance du 23 novembre 2023, « Reconnaissance d'entités nommées ») dans [Voyant Tools](#) depuis votre ordinateur ;

Voyant Tools est un environnement en ligne de lecture et d'analyse de textes numériques.
Traduction française d'Aurélien Berra

2. Combien de *tokens* et de *types* ce corpus contient-il ?

Résumé Documents Syntagmes ?

Ce corpus contient 7 documents, 11,106 mots et 2,668 formes verbales uniques. Il a été créé maintenant.

Taille du document :

- Ordre décroissant : Le_Tour_du_monde_en_quatr... (2410); Le_Tour_du_monde_en_quatr... (1824); Le_Tour_du_monde_en_quatr... (1740); Le_Tour_du_monde_en_quatr... (1509); Le_Tour_du_monde_en_quatr... (1277)
- Ordre croissant : Le_Tour_du_monde_en_quatr... (1069); Le_Tour_du_monde_en_quatr... (1277); Le_Tour_du_monde_en_quatr... (1277); Le_Tour_du_monde_en_quatr... (1509); Le_Tour_du_monde_en_quatr... (1509)

Élément

Tokens : 11 106

Types : 2 668

3. Quelle est la fréquence absolue du mot **fogg**? 107

Cirrus

Termes

Liens

?

			Terme	Total	Tendance
<div>+</div>	<div>☑</div>	1	fogg	107	<div></div>
<div>+</div>	<div>☐</div>	2	phileas	71	<div></div>
<div>+</div>	<div>☐</div>	3	heures	48	<div></div>
<div>+</div>	<div>☐</div>	4	passpartout	44	<div></div>
<div>+</div>	<div>☐</div>	5	répondit	43	<div></div>
<div>+</div>	<div>☐</div>	6	monsieur	37	<div></div>
<div>+</div>	<div>☐</div>	7	qu'il	32	<div></div>
<div>+</div>	<div>☐</div>	8	mr	31	<div></div>
<div>+</div>	<div>☐</div>	9	jours	30	<div></div>
<div>+</div>	<div>☐</div>	10	licence	28	<div></div>

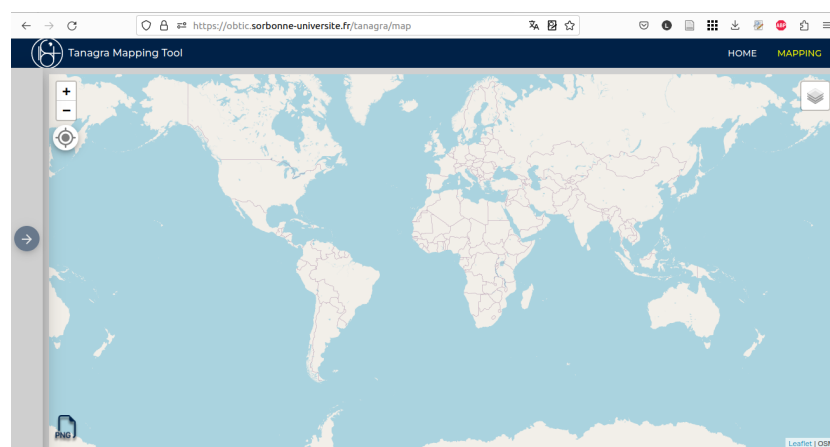
?

2,456

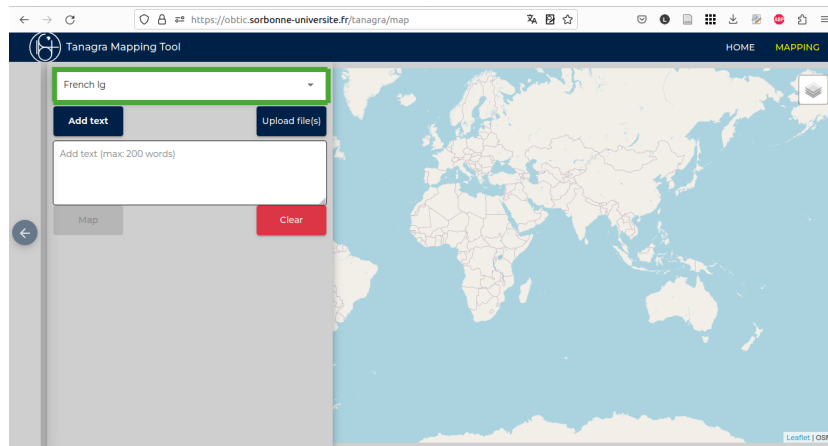
Question 10

(3 points)

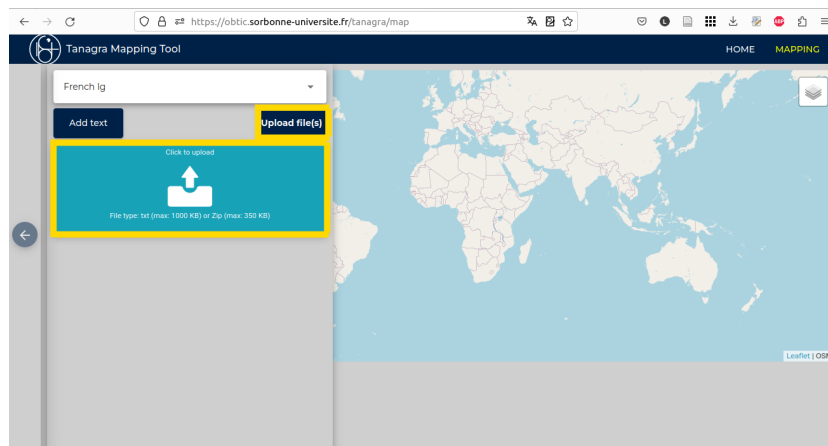
1. Dirigez-vous vers la page de l'outil [Tanagra](https://obtic.sorbonne-universite.fr/tanagra/map).



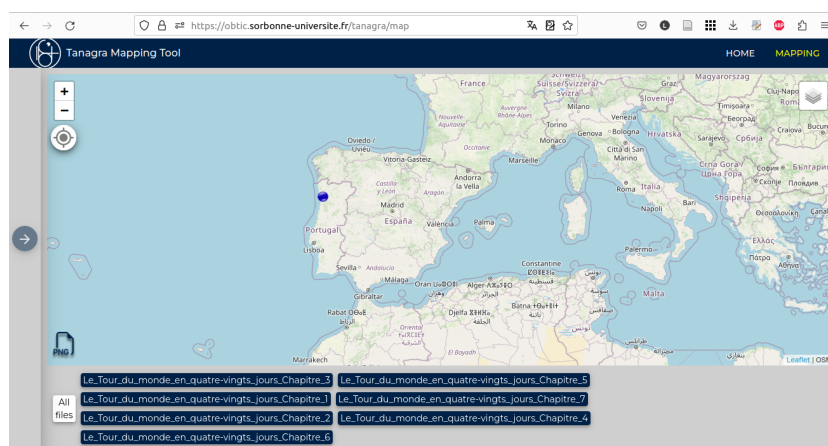
2. Sélectionnez **French_lg** comme modèle de reconnaissance d'entités nommées :



et chargez le même corpus que dans la question 9.



3. Parmi les entités nommées récupérées par Tanagra, en trouvez-vous une localisée en Espagne ? Si oui, laquelle ?



Non, aucune entité nommée localisée en Espagne n'a été récupérée.