



RECONNAISSANCE OPTIQUE DE CARACTÈRES DANS LES PERMIS CITES À L'AIDE DE L'APPRENTISSAGE MACHINE

22-23 février 2023

BUREAU DU DIRIGEANT PRINCIPAL DES DONNÉES
DIRECTION GÉNÉRALE DE LA POLITIQUE STRATÉGIQUE

PAR :

ARPIT RATHORE – Scientifique des données

LAKSHAY GOEL – Développeur de logiciels (stagiaire)

MICHAEL HOCKEY – Développeur de logiciels (stagiaire)



BUREAU DU DIRIGEANT PRINCIPAL DES DONNÉES

QU'EST-CE QUE LA ROC?

- L'acronyme « ROC » désigne la **reconnaissance optique de caractères**. C'est une technologie qui reconnaît le texte dans une image numérique.
- La ROC **convertit les documents imprimés en documents-texte lisibles par une machine**.



BUREAU DU DIRIGEANT PRINCIPAL DES DONNÉES

QUE SONT LES PERMIS CITES?

- La Convention sur le commerce international des espèces de faune et de flore sauvages menacées d'extinction (CITES) est **une entente dont l'objectif est de veiller à ce que le commerce international des animaux et des plantes sauvages ne menace pas leur survie.**
- Un permis CITES est obligatoire pour l'importation et l'exportation de produits dérivés des espèces protégées par la CITES. Cette protection s'applique aux espèces inscrites à la CITES, quelle que soit la forme de l'individu (mort ou vivant).
- Le Service canadien de la faune (SCF) d'Environnement et Changement climatique Canada supervise l'administration des permis CITES pour le Canada et les autres pays importateurs ou exportateurs effectuant des échanges commerciaux au Canada.



BUREAU DU DIRIGEANT PRINCIPAL DES DONNÉES

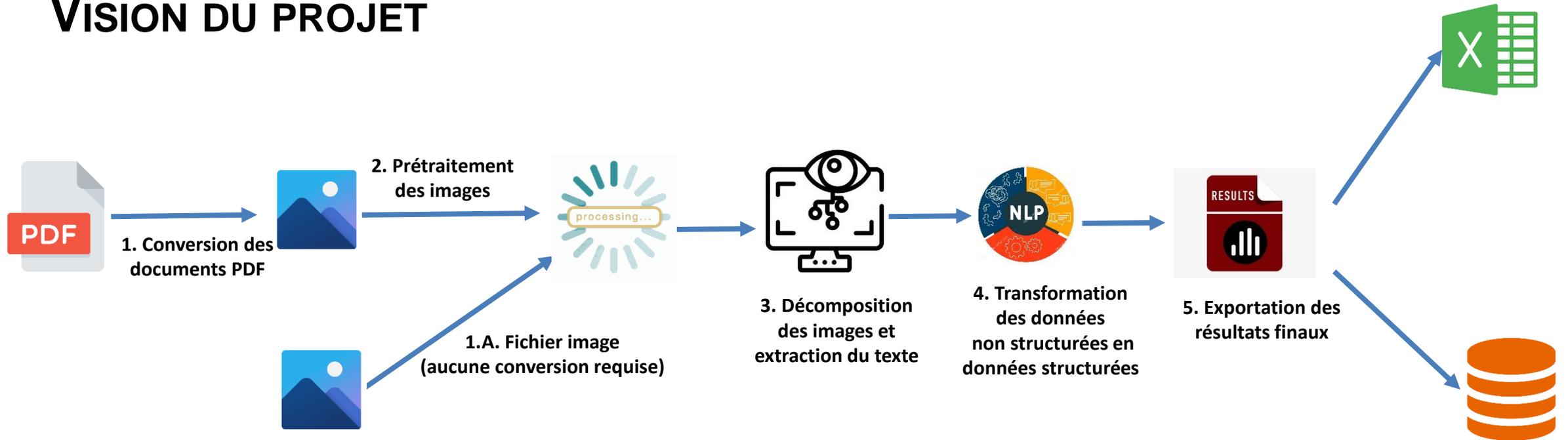
VISION DU PROJET

- Soutenir la capacité du SCF d'accéder aux données des permis CITES.
- Convertir, au moyen d'une technologie de vision informatique, les milliers de permis CITES numérisés de multiples années et de divers pays pour former une base de données numériques.
- Examiner et évaluer des options de technologie de vision informatique (Azure Form Recognizer et des bibliothèques Python) pour choisir la solution optimale.
- Utiliser les données converties des permis CITES à des fins d'analyse approfondie.
- Mettre au point un mécanisme pour que le SCF puisse téléverser de futurs permis dans la base de données en les transformant.



BUREAU DU DIRIGEANT PRINCIPAL DES DONNÉES

VISION DU PROJET



BUREAU DU DIRIGEANT PRINCIPAL DES DONNÉES

JALONS DU PROJET

1. Analyse des options

- L'outil Microsoft Azure Form Recognizer a été mis à l'essai pour la ROC, mais il ne convenait pas à la nature des permis CITES, car il a besoin d'au moins quatre échantillons de données pour procéder à un nouvel entraînement avant de passer à l'extraction.
- Des modèles Python préentraînés pour la ROC, comme PyTesseract et EasyOCR, ont été mis à l'essai, mais l'extraction était brouillonne et la transformation en données structurées était impossible.
- Il a donc été déterminé que la solution optimale était une conception sur mesure en Python. Le flux de projet et la conception des outils personnalisés intègrent des bibliothèques Python comme OpenCV, EasyOCR, Pillow, pdf2image, NLTK, etc.

2. Conversion des fichiers

- À l'heure actuelle, les permis CITES sont numérisés et conservés sous forme de documents PDF de plusieurs pages.
- Nous avons créé un module de validation des données pour vérifier le format du fichier entrant.
- Nous convertissons les documents PDF en images en les envoyant au module de conversion des données; sinon, les fichiers sont envoyés au module de prétraitement d'image.
- Les principaux ensembles Python servant à cette tâche sont pdf2image, PIL, etc.



BUREAU DU DIRIGEANT PRINCIPAL DES DONNÉES

VISION DU PROJET

3. Prétraitement d'image

- Les images converties passent par le module de prétraitement d'image, qui crée des images uniformes, corrige l'obliquité, élimine le bruit, convertit le RVB en monochrome, dilate les images, etc. [[VOIR LA DIAPO 9](#), image 1 à image 2]
- Ce prétraitement aide à créer des données normalisées pour le module de détection de modèles.
- Les principaux ensembles Python servant à cette tâche sont OpenCV, NumPy, etc.

4. Détection de modèles, décomposition d'image et extraction de texte

- Les images prétraitées passent par le module de détection de modèles de permis CITES, qui facilite la décomposition des images : le module de décomposition se base sur les modèles détectés pour diviser l'image en n images de plus petites tailles. [[VOIR LA DIAPO 9](#), image 3]
- Les n images décomposées subissent un processus d'extraction de texte par ROC, qui va chercher tout le texte présent dans chacune des petites images. [[VOIR LA DIAPO 10](#), image 4 à image 5]
- Les principaux ensembles Python servant à cette tâche sont OpenCV, Pandas, PyTesseract, EasyOCR, etc.



BUREAU DU DIRIGEANT PRINCIPAL DES DONNÉES

VISION DU PROJET

5 et 6. Transformation des données brutes en données structurées et exportation des données structurées (en cours de planification)

- Des méthodes de traitement du langage naturel seront utilisées pour convertir les données brutes en données structurées, afin de gérer les fautes d'orthographe et les incohérences dans les données entrées.
- Les principaux ensembles Python servant à cette tâche sont NLTK, NumPy, Pandas, RegEx/Regular Expressions et NER, etc.



BUREAU DU DIRIGEANT PRINCIPAL DES DONNÉES

This is the original document image, showing a form titled 'CONVENTION ON INTERNATIONAL TRADE IN ENDANGERED SPECIES OF WILD FAUNA AND FLORA'. The form is filled with text and includes a red circular stamp at the bottom right. The image is slightly blurred and has a white background.

1001878

1. Image d'origine

This is the pre-processed document image, where the form is clearer and the text is more legible. The red circular stamp is still present. The background is white, and the text is black.

1001878

2. Image prétraitée

This is the document image with green bounding boxes overlaid on it, indicating the detection of models and decomposition. The bounding boxes are rectangular and cover various sections of the form. The red circular stamp is still present.

1001878

3. Détection de modèles et décomposition



BUREAU DU DIRIGEANT PRINCIPAL DES DONNÉES

	A	B	C	D	E	F	G	H	I	J
0.jpg	11. Export	IDELVAUX	[Belgique]							
1.jpg	No / No									
2.jpg	9. Masse nette (kg)	Net mass (kg)								
3.jpg	10. Quantité	Quantity	1 pce							
4.jpg	11. Annxc CITES	CITES Appendix								
5.jpg	12. Annxc CE EU	Annex								
6.jpg	13. Origine	Source								
7.jpg	14. Objet	Purpose								
8.jpg	15. Pays d'origine	1 Country of origin	États-Unis							
9.jpg	16. N° de permis	Permi	No							
10.jpg	17. Date de délivrance	Date of issue								
11.jpg	18 Pays de dernière réexportation	'Country' of last re-export								
12.jpg	2. Dernier jour de validité	Last day of validity								
13.jpg	19. No du certificat	Certificate No								
14.jpg	22o. Date de délivrance	Date of issue								
15.jpg	k1. Nom scientifique de Y'espèce	Scientific name of species								
16.jpg	22. Nom commun de Tespèce	Common name of species								

4. Images découpées

5. Données extraites brutes



DES QUESTIONS?

Veillez communiquer avec moi :

Arpit Rathore

Scientifique des données

Bureau du dirigeant principal des données

Direction générale de la politique stratégique, Environnement et Changement climatique Canada

Arpit.Rathore@ec.gc.ca
