## L'IA pour le classement automatique des articles collectés via web scraping – cas de la VSI de la Plateforme ESV



« La Plateforme est en charge d'assurer une veille internationale sur les dangers sanitaires susceptibles d'avoir un impact sanitaire et/ou économique » (extrait de la convention cadre).



Isabelle Pieretti (CIRAD) et Antoine Marullaz (INRAE)
Journée Vectopole Sud, Montpellier, 18-09-2025

https://plateforme-esv.fr/

<u>linkedin</u> in















## Objectifs

Anticiper et mettre à jour des connaissances pour le gestionnaire du risque et les acteurs de la surveillance sur le territoire métropolitain et dans les DROMs.

- 1) Augmenter le nombre d'Organismes Nuisibles ciblés (couverture optimale de la VSI)
- 2) Produire une veille adaptée et la valoriser

### Comment?

Utiliser et développer des outils basés sur l'IA et la fouille de texte pour améliorer en continu le pipeline de la VSI (collaborations VSI avec SI, autres informaticiens, projets de recherche connexes).

#### Pour:

- 1. Réduire le nombre d'articles « non pertinents » (hors scope VSI)
- 2. Automatiser les tâches qui peuvent l'être
- 3. Optimiser les requêtes et les sources



ectopole Sud, Montpellier, 18-09-2025

## Productions (exemple)





BHV-SV 2025/17 Mois d'Avril 2025 semaine 17 du 21 au 27 avril 2025

### Bulletin Hebdomadaire de Veille Sanitaire Internationale en Santé Végétale

Le Bulletin Hebdomadaire de Veille sanitaire internationale en Santé Végétale (BHV-SV) s'inscrit dans l'activité de veille sanitaire internationale menée dans le cadre de la Plateforme ESV (Plateforme d'Épidémiosurveillance en Santé Végétale -https://www.plateforme-esv.fr/). Le BHV-SV sélectionne et résume des actualités sanitaires et scientifiques en santé végétale qui sont parues dans la semaine.

ATTENTION : Le contenu du document n'engage pas les membres de la Plateforme ESV.



Attribution - Pas d'Utilisation Commerciale CC BY-NC-ND

Code juridique

Conformément aux productions réalisées par la Plateforme d'Épidémiosurveillance en Santé Végétale (ESV), celle-ci donne son droit d'accès à une utilisation partielle ou entière par les médias, à condition de ne pas apporter de modification, de respecter un cadre d'usage bienveillant et de mentionner la source © https://plateforme-esv.fr/

Confiance + est un indicateur sur la crédibilité des sites de diffusion (+ : peu fiable à +++ : très fiable, source officielle majoritairement).

#### **Sommaire**

Veille non ciblée	2
Agrilus planipennis	3
Agrilus anxius	3
Xylella fastidiosa	4

#### Popillia japonica

#### Veille sanitaire

En Suisse, dans le canton de Schwytz, différentes mesures dans la zone d'infestation et dans la zone tampon ont été prises via un arrêté général depuis la découverte de *Popillia japonica* au cours de l'été 2024.

Titre	Categorie	PaysSujet	Fiabilite	Lien
Mesures contre les scarabées japonais dans la région de Sägel	Réglementation,Evolution de l'état sanitaire	Suisse	Officielle	lien

résumé de l'article

autres informations

#### Veille scientifique

Voici un rapport de synthèse sur la biologie, la propagation, le potentiel de nuisance, les mesures de surveillance et de lutte concernant *Popillia japonica* (langues EN, FR, IT).

Titre	Categorie	Lien
Japanese beetle: current information on its biology, legal	Synthèse et	lien
bases and control measures	sensibilisation	

#### Ceratocystis platani

#### Veille sanitaire

Plante & Cité a publié, le 18 mars 2025, une version enrichie du guide destiné à aider les gestionnaires et propriétaires de platanes à prévenir et limiter la dissémination la maladie due au champignon Ceratocystis platani.

Titre	Categorie	PaysSujet	Fiabilite	Lien
Un kit pour aider à lutter contre le chancre coloré du platane	Méthode et mesure de lutte,Méthode, outil et mesure de surveillance	France	Agronomique	lien



Vectopole Sud, Montpellier, 18-09-2025

## Productions (exemple)

#### Spodoptera frugiperda

Popillia japonica

Candidatus Liberibacter et ses vecteurs

Tourneyella parvicornis

Xylella fastidiosa

#### **BULLETIN MENSUEL N°68**

Plateforme ESV

Mars 2025



Le bulletin d'Épidémiosurveillance en Santé Végétale est une revue des actualités concernant la santé du végétal en Europe et à l'International. Il contribue à faciliter l'accès aux informations concernant la santé des végétaux et leur diffusion. Le bulletin est validé au préalable par une cellule éditoriale composée d'experts scientifiques et de collaborateurs partenaires ayant un rôle de conseillers.



Attribution - Pas d'Utilisation Commerciale - Pas de modification CC BY-NC-ND

00 01 110 111

Code juridique

Conformément aux productions réalisées par la Plateforme d'Épidémiosurveillance en Santé Végétale (ESV), celle-ci donne son droit d'accès à une utilisation partielle ou entière par les médias, à condition de ne pas apporter de modification, de respecter un cadre d'usage bienveillant et de mentionner la source © https://www.plateforme-esv.fr/

Le bulletin est enregistrable dans son format html sur votre ordinateur. Toutes les cartes du bulletins sont interactives (zoom/ dézoom, informations cliquables directement sur la carte).

#### Sommaire

bollillane			
Sujet phytosanitaire	Zone géographique	Cultures	Nature de l'information
Xylella fastidiosa	Espagne	Multi- espèces	Évolution de l'état sanitaire et réglementaire
Popillia japonica	Italie	Multi- espèces	Évolution de l'état sanitaire et réglementaire
Spodoptera frugiperda	Grèce	Multi- espèces	Évolution de l'état sanitaire
Tourneyella parvicornis	Italie	Pins	Évolution de l'état sanitaire et réglementaire
Tourneyella parvicornis	Italie	Pins	Outil de détection des symptômes de dépérissement forestier causés par la cochenille tortue du pin
Candidatus Liberibacter spp. et ses vecteurs	Monde	Agrumes	Gestion collective d'une maladie végétale incurable : cas du Huanglongbing
Candidatus Liberibacter spp. et ses vecteurs	Portugal	Agrumes	Estimation du risque épidémiologique

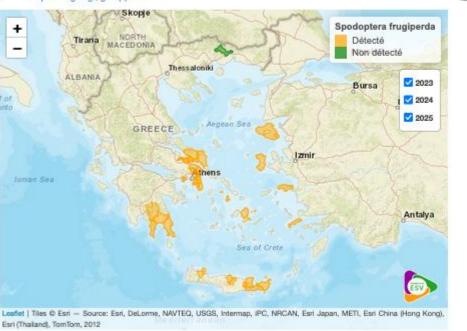
### Spodoptera frugiperda

Grèce / Multi-espèces / Évolution de l'état sanitaire

La présence de Spodoptera frugiperda en Europe a été établie pour la première fois en 2021 dans toutes les îles espagnoles de l'archipel des Canaries (Lanzarote, Fuerteventura, Gran Canaria, Tenerife, La Palma, El Hierro et La Gomera, source OEPP), puis en 2023 en Roumanie (comté de Calarasi, source OEPP), à Chypre (districts de Limassol et Larnaca, source OEPP), sur l'île portugaise de Madère (OEPP) et en Grèce. En Grèce, l'organisme de quarantaine prioritaire pour l'UE a été découvert dans le cadre de la surveillance officielle via un dispositif de piégeage à phéromones. Le ravageur a ainsi été capturé dans diverses préfectures du pays. En Laconie et Attique orientale pour la zone continentale, et dans les îles suivantes: Crète (Héraklion, Chania ou La Canée, Lassithi), Eubée, Lesbos, Salamine, Kos, Chios, Samos, Naxos et Syros. Dans la préfecture de Xanthi (nord de la Grèce), le réseau de pièges a été établi dans des cultures de mais, de riz et d'autres céréales (blé, avoine, orge) ainsi que des cultures de luzerne. Le ravageur n'a pas été découvert à ce jour à Xanthi.

texte

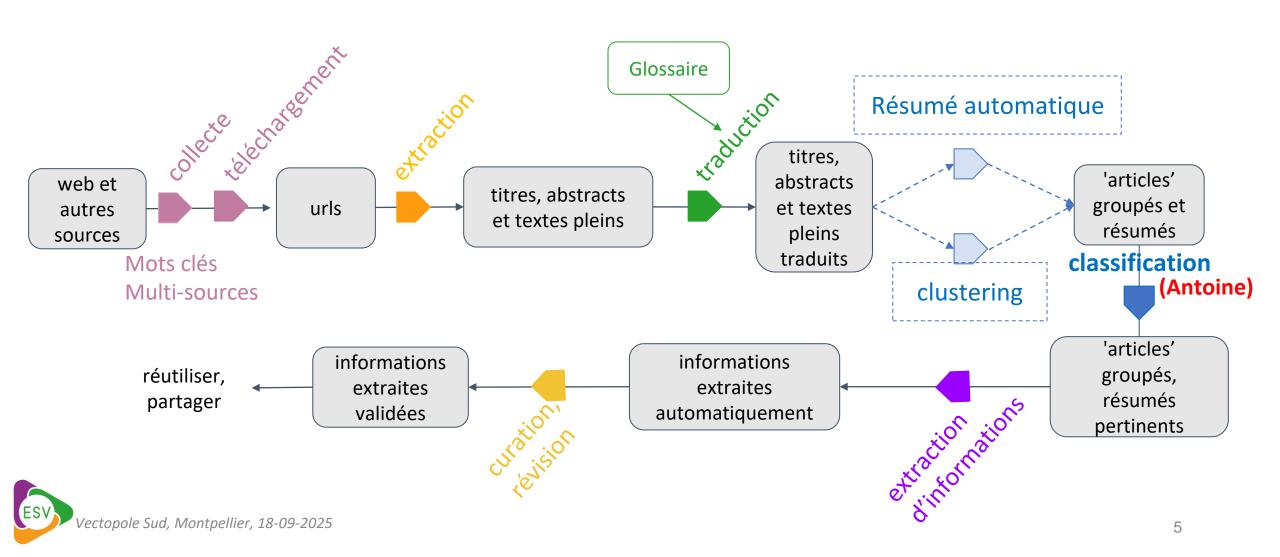
Sources: pamth.gov.gr, gd.eppo.int.



carte interactive

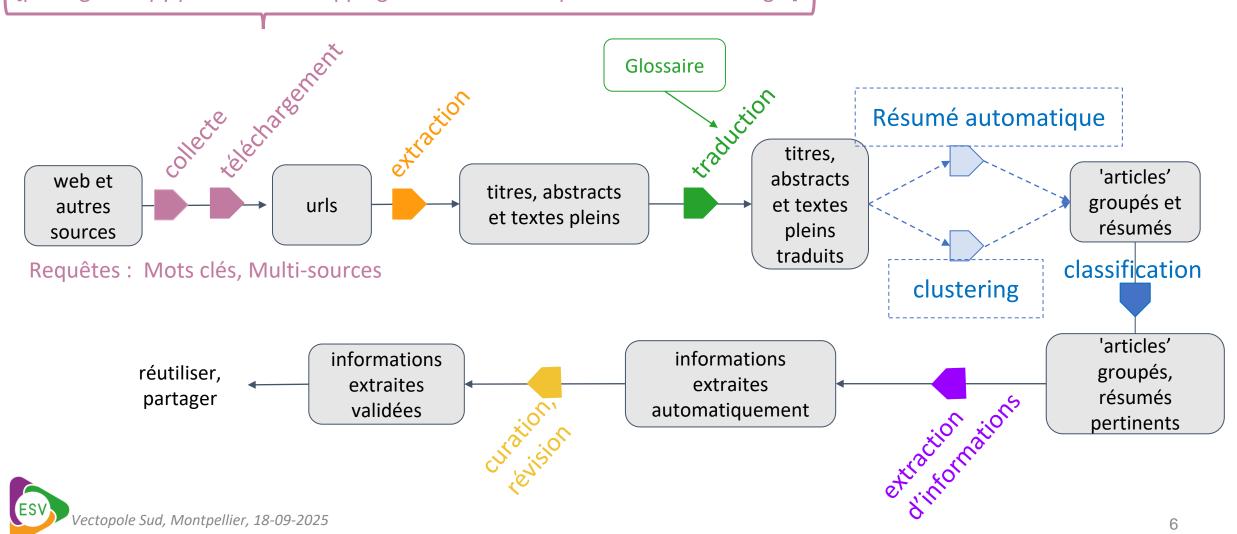
Figure 1 : Carte des détections de Spodoptera frugiperda dans l'union européene, à noter que les zones notées en "Détecté" pour 2025 sont issues d'un document de 2025 dans lequel les dates de détection ne sont pas précisées. Sources pour les zones non relayées dans les BM précédents : pamth.gov.gr, gd.eppo.int.

### PIPELINE VSI —> Lancé tous les lundis matins

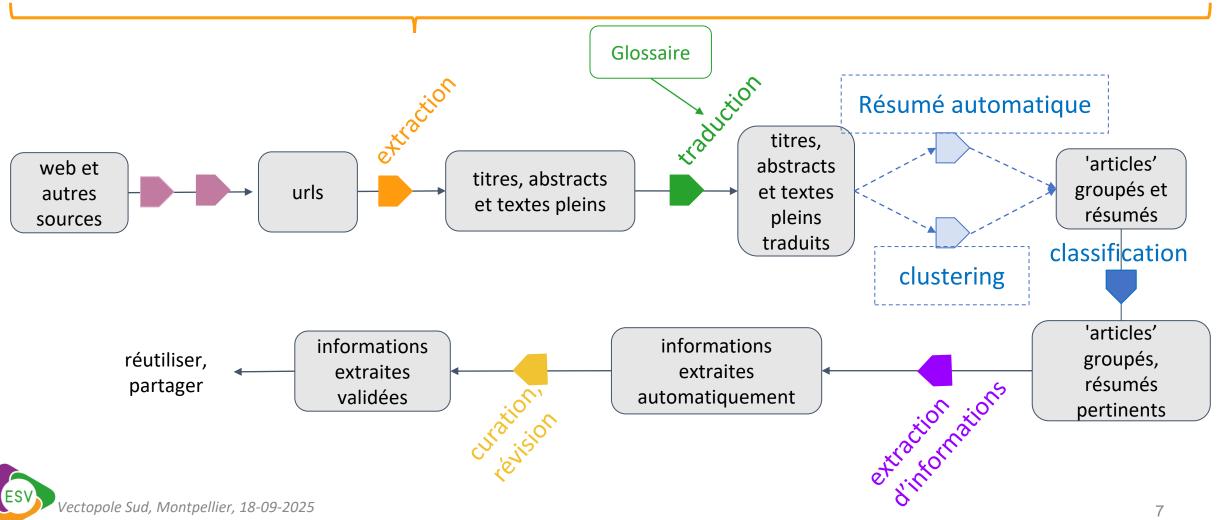


### Collecte, téléchargement automatique des URLs avec scripts python

[package Scrapy pour le webscrapping et API scaleSERP pour les sources Google]



Extraction document via script python utilisant librairie Trafilatura + extraction contenu page web au format XML\_TEI



### **Étape 1 : OUVERTURE d'UNE URL** (exemple)



https://www.vivicastellanagrotte.it/index.php/notizietop/16838monitoraggio-xylella-tre-piante-positive-a-castellana-grotte

HOME

WEBTV

DAL PALAZZO

**CULTURA ~** 

**SPORT** GROTTE DI CASTELLANA > LA CITTADINA ~

NEWSLETTER

CONTATTI



Apnée du sommeil : un oreiller ingénieux fait fureur au France

NaturallyYou

**OUVRIR** 



å COMUNICATO STAMPA ▷ NOTIZIE 🖰 09 SETTEMBRE 2022

#### Monitoraggio xylella - Tre piante positive a Castellana-Grotte

Riceviamo comunicazione del rinvenimento di tre piante di ulivo positive nella periferia nord-est di Castellana-Grotte. Dopo Locorotondo, Monopoli, Polignano a mare e Alberobello - dunque - la Città

delle Grotte è il quinto Comune della Città Metropolitana di Bari nel quale sia stata rilevata la presenza del batterio xylella fastidiosa. Le tre piante, tra loro adiacenti, sono al confine tra la "Zona cuscinetto Polignano a mare" e la zona indenne. Il primo aggiornamento dopo la pausa estiva comprende anche altre 13 piante infette, due delle quali nella (ormai ex) zona indenne di Polignano a mare, altre 6 in zona contenimento (5 a Fasano e 1 a Martina Franca); le restanti cinque son in zona cuscinetto a Polignano a mare (3), a Monopoli (1) e una ad Alberobello. In quest'ultimo caso si tratta del primo ritrovamento a ovest del centro abitato della capitale dei trulli, a circa 1 Km in direzione Noci. Le analisi, effettuate dal Centro di ricerca sperimentazione e formazione in agricoltura "Basile Caramia" di Locorotondo su campioni rilevati il 23 agosto 2022, sono datate al 29 agosto u.s.



































# Étape 2 : OUVERTURE code source en HTML de la page web



HOME WEBTV CULTURA ~ SPORT NEWSLETTER CONTATTI DAL PALAZZO GROTTE DI CASTELLANA V LA CITTADINA ~ <h2 itemprop="name"> 130 Monitoraggio xylella - Tre piante positive a Castellana-Grotte 131 </h2> 132 </div> 133 134 135 138 137 138 <div itemprop="articleBody"> <div class="17ghb35v kjdc1dyq kmwttqpk gh25dzvf jikcssrz n3t5jt4f">Riceviamo comunicazione del rinvenimento di tre piante di ulivo positive nella periferia nord-est di Caste. 140 <a class="a2a\_button\_facebook"></a> <a class="a2a button twitter"></a> 142 <a class="a2a button telegram"></a> 143 <a class="a2a\_button\_whatsapp"></a> 144 <a class="a2a\_button\_threads"></a> 145 <a class="a2a button email"></a> 148 <a class="a2a button print"></a> 147 </span> </div> </div> 148 149 150 151 152 class="previous"> 153 <a class="hasTooltip" title="I.I.S.S. ''dell'Erba'' - A.S. 2022-2023: si comincia" aria-label="Articolo precedente: I.I.S.S. ''dell'Erba'' - A.S. 2022-2023: si comincia" hret 154 <span class="icon-chevron-left" aria-hidden="true"></span aria-hidden="true">Indietro</span> 155 









### Étape 3 : EXTRACTION CONTENU TEXTUEL au format XML\_TEI via Trafilatura



### Monitoraggio xylella - Tre piante positive a Castellana-

#### Grotte

Riceviamo comunicazione del rinvenimento di tre piante di ulivo positive nella periferia nord-est di Castellana-Grotte. Dopo Locorotondo, Monopoli, Polignano a mare e Alberobello - dunque - la Città delle Grotte è il quinto Comune della Città Metropolitana di Bari nel quale sia stata rilevata la presenza del batterio xylella fastidiosa. Le tre piante, tra loro adiacenti, sono al confine tra la "Zona cuscinetto Polignano a mare" e la zona indenne. Il primo aggiornamento dopo la pausa estiva comprende anche altre 13 piante infette, due delle quali nella (ormai ex) zona indenne di Polignano a mare, altre 6 in zona contenimento (5 a Fasano e 1 a Martina Franca); le restanti cinque son in zona cuscinetto a Polignano a mare (3), a Monopoli (1) e una ad Alberobello. In quest'ultimo caso si tratta del primo ritrovamento a ovest del centro abitato della capitale dei trulli, a circa 1 Km in direzione Noci. Le analisi, effettuate dal Centro di ricerca sperimentazione e formazione in agricoltura "Basile Caramia" di Locorotondo su campioni rilevati il 23 agosto 2022, sono datate al 29 agosto u.s.









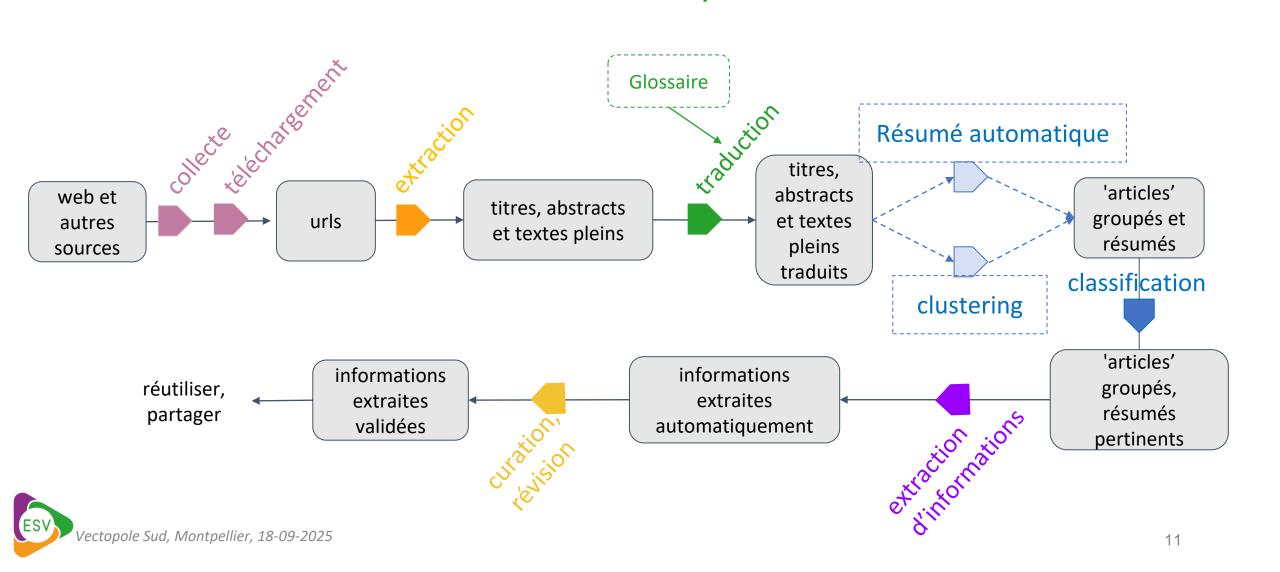




Avanti 🕽

"<TEI xmlns=\"http://www.tei-c.org/ns/1.0\">\n <teiHeader>\n <fileDesc>\n <title type=\"main\">Monitoraggio xylella - Tre piante positive a Castellana-Grotte</title>\n <author>Comunicato Stampa</author>\n </titleStmt>\n <publicationStmt>\n \n </publicationStmt>\n <notesStmt>\n <note type=\"fingerprint\">ASwdniSgGiFYGI51Y9b3Kxau9AM=</note>\n </notesStmt>\n <sourceDesc>\n <bib|>Monitoraggio xylella - Tre piante positive a Castellana-Grotte. vivicastellanagrotte.it, 2022-09-09</bibl>\n <bibl type=\"sigle\">vivicastellanagrotte.it, 2022-09-09</bibl>\n <biblFull>\n <title type=\"main\">Monitoraggio xylella - Tre piante positive a Castellana-Grotte</title>\n <a href="mailto:catellana-Grotte">a castellana-Grotte</title>\n <a href="mailto:catellana-Grotte">a castellana-Grotte</a>/title>\n <a href="mailto:catellana-Grotte">a castellana-Grotte</a>/ <publicationStmt>\n <publisher>vivicastellanagrotte.it (vivicastellanagrotte.it)</publisher>\n <ptr type=\"URL\" target=\"http://www.vivicastellanagrotte.it/index.php/notizietop/16838-monitoraggio-xylella-tre-piante-positive-a-castellana-<abstract>\n Riceviamo comunicazione del rinvenimento di tre piante di ulivo positive nella periferia nord-est di Castellana-Grotte. Dopo Locorotondo, Monopoli, Poligna...\n </abstract>\n <textClass>\n <keywords>\n <term type=\"tags\">webtv, web, tv, castellana, grotte, fanove, madonna della vetrana, puglia, bari</term>\n </keywords>\n </textClass>\n </profileDesc>\n <encodingDesc>\n <applifo>\n <application version=\"1.3.0\" ident=\"Trafilatura\">\n </encodingDesc>\n </teiHeader>\n <text>\n <body>\n <div type=\"entry\">Riceviamo comunicazione del rinvenimento di tre piante di ulivo positive nella periferia nord-est di Castellana-Grotte. Dopo Locorotondo, Monopoli, Polignano a mare e Alberobello - dunque - la Città delle Grotte è il quinto Comune della Città Metropolitana di Bari nel quale sia stata rilevata la presenza del batterio xylella fastidiosa.<hi rend=\"#i\">Le tre piante, tra loro adiacenti, sono al confine tra la "Zona cuscinetto Polignano a mare" e la zona indenne. Il primo aggiornamento dopo la pausa estiva comprende anche altre 13 piante infette, due delle quali nella (ormai ex) zona indenne di Polignano a mare, altre 6 in zona contenimento (5 a Fasano e 1 a Martina Franca); le restanti cinque son in zona cuscinetto a Polignano a mare (3), a Monopoli (1) e una ad Alberobello. In quest'ultimo caso si tratta del primo ritrovamento a ovest del centro abitato della capitale dei trulli, a circa 1 Km in direzione Noci.</hi>Le analisi, effettuate dal Centro di ricerca sperimentazione e formazione in agricoltura "Basile Caramia" di Locorotondo su campioni rilevati il 23 agosto 2022, sono datate al 29 agosto u.s.</div>\n <div type=\"comments\"/>\n  $</body>\n </text>\n </TEI>"$ 

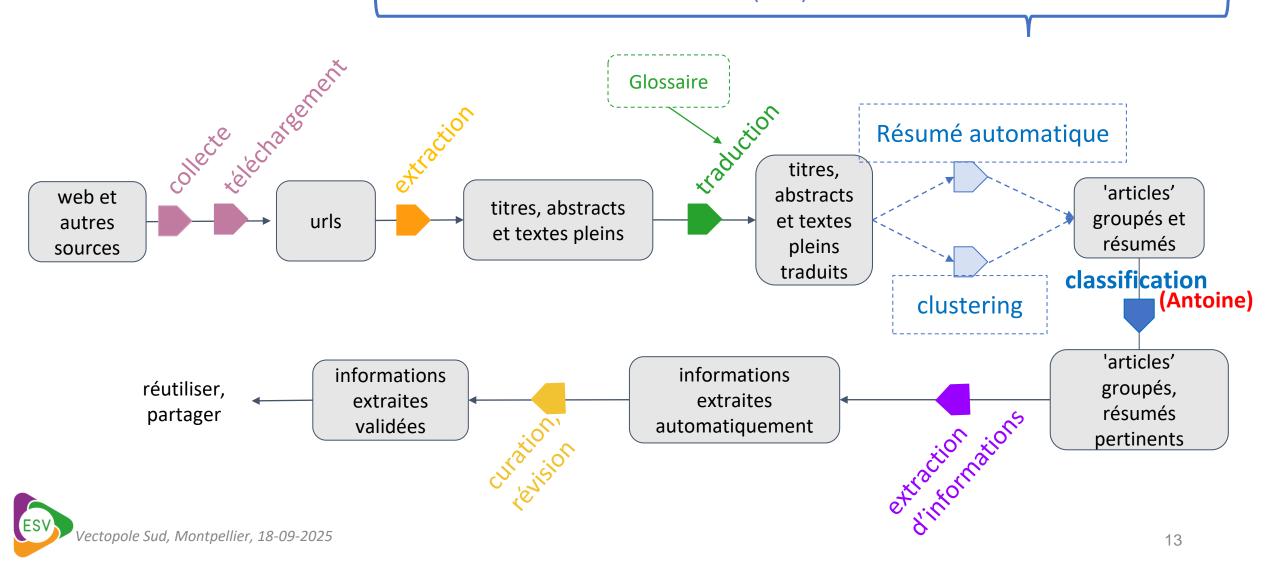
RENDRE L'INFORMATION COMPREHENSIBLE : traduction automatique du titre et du texte plein (EN) via script python qui appelle API Google Translate



https://www.alaftharia		Titre publié (angla	-	Titre collecté (multi-langues)	Occio		
https://www.eleftheria.	.gr/peth true	Les platanes meurent dans la vallée	ае тетірі	Πεθαίνουν τα πλατάνια στην Κοιλάδα των Τεμπών - Ελει	grec grec		
https://33.fsvps.gov.ru	/news/ true	Under the supervision of the Office	of Rosselkhoznadzor in Vladimir, measures have been taken to localize and eli	Под контролем Управления Россельхознадзора во	russe		
https://rezeknesnovad	s.lv/osu true	Surveillance de l'agrile du frêne dan	is la région de Rēzekne	Ošu smaragdzaļās krāšņvaboles monitorings letton			
https://rosleshoz.gov.r	u/news/ true	Emerald ash borer found in Krasnoa	rmeyskoye forestry in Saratov region	В Красноармейском лесничестве Саратовской области			
https://www.gov.pl/we	b/wiori true	Control activities within the EABRAC	E project - monitoring of the emerald ash borer in the Podkarpackie Voivodesh	Działania kontrolne w ramach projektu EABRACE p	olonais		
https://www.gov.me/cl	anak/la true	Analysis of grapevine leaf samples of	did not reveal the presence of the phytoplasma Flavescence dorée	Analizom uzoraka listova vinove loze nije utvrđeno	serbe		
https://dentrolanotiziat	oreak.it/ true	Agriculture alert: small, but devastat	ting, beetle invasion in Piedmont - Inside the news	è piccolo, ma devasta il Piemonte, l'invasione del coleotter	italien		
https://www.juntadean	dalucia true	First detections of Popillia japonica i	in Galicia and France - RAIF	Primeras detecciones de Popillia japonica en Galicia y	espagnol		
https://www.yvorne.ch	<u>/</u> true	Bienvenue sur le site Internet officie		Commune d'Yvorne: Bienvenue sur le site Internet officiel	français		
	33 (2025) : 963		Citrus Forum on pests and technology	Panamá será sede del IV Foro Internacional de Cítricos			
capture image	e OVIDE (IHM)/	comité éditorial	es for the control of HLB in citrus plants was established.	Se estableció un nuevo plan integral de medidas para			
https://www.lagaceta.c	com.ar/ true	HLB: Pilcomayo regained its vector-	free status	HLB: Pilcomayo recuperó el status de área libre del vector			
https://www.fecier.org	.ar/notic true	The citrus leafhopper, carrier of the	dreaded HLB bacteria, is already roaming the orange groves of San Pedro and	Federacion del Citrus de Entre Rios			
https://papers.ssrn.com	m/sol3/ true	Identification of Soil Microbial Comm	nunities in Huanglongbing Infected Mandarin Orange Orchards Through Metag	Identification of Soil Microbial Communities in anglais			
https://www.frontiersin	n.org/jo true	Frontiers   The Application of Bacillu	us amyloliquefaciens and Arbuscular Mycorrhizal Fungi Displays Curative Effect	The Application of Bacillus amyloliquefaciens and			
https://qdxinshuoyuan	.com/3 true	General Administration of Customs	Announcement No. 165 of 2025 (Announcement on Plant Quarantine Requirem	海关总署公告2025年第165号(关于中国柚子出口新西兰植物)	chinois		
https://www.sciencedi	rect.co true	Explainable ensemble learning for pr	redicting pine wilt disease spread	Explainable ensemble learning for predicting pine wilt			
https://www.haber236	.com/ya true	Study on leaf blight disease		Yaprak yanıklığı hastalığına yönelik çalışma turo			

#### TRIER-MODIFIER-CLASSER LES INFOS PERTINENTES:

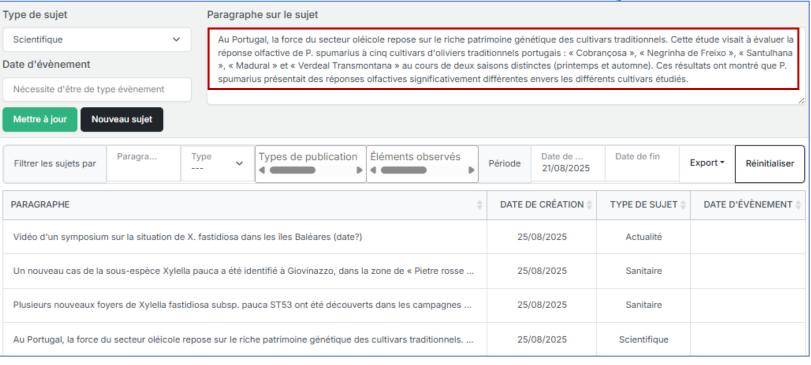
- Classification automatique (présentation Antoine)
- Ecriture d'un résumé automatique du texte plein (à tester en routine)
- Clustering automatique des doublons (à tester en routine).
- Modification via interface OVIDE (IHM)



 résumé automatique via DistilBART (model encodeur-décodeur) (à tester en routine)

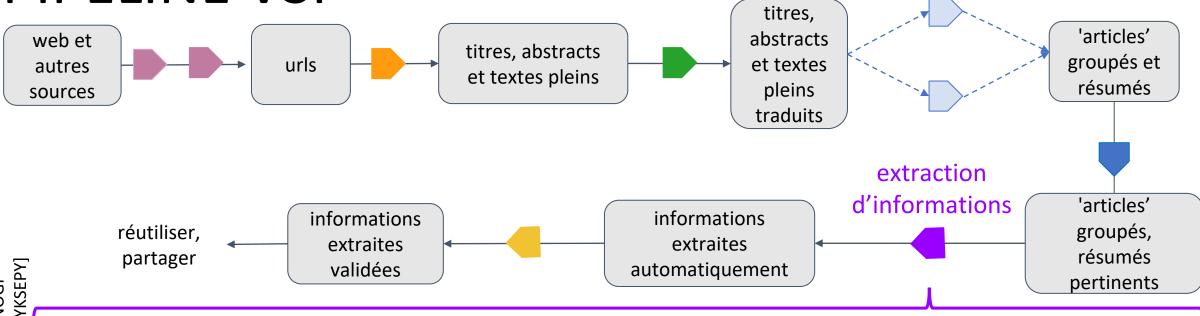
 clustering automatique des doublons via algo DBSCAN (à tester en routine)

### Interface IHM OVIDE



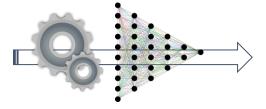
ID \$	URL \$	ENVOYÉ AU CE 🌲	TITRE PUBLIÉ	
254509	https://www.letemps.ch/scienc	true	Terreur des vignobles, le scarabée japonais débarque en Suisse romande	
254507	https://www.rts.ch/info/regions/	true	Un premier foyer de scarabée japonais identifié à Genève	
253966	https://www.watson.ch/fr/suiss	false	Un foyer de cet insecte nuisible identifié à Genève	
253964	https://www.tdg.ch/geneve-un	false	Un premier foyer de cet insecte invasif identifié dans le canton de Genève	
253960	https://www.letemps.ch/scienc	false	Terreur des vignobles, le scarabée japonais débarque en Suisse romande - Le Temps	





Campagne d'Annotation par des experts\* pour assister la VSI par des méthodes d'Intelligence Artificielle, notamment via des outils issus de la recherche en traitement automatique de la langue (TAL) 
Corpus Epidemiomonitoring Of Plant (EPOP) [Dataverse Plateforme ESV] 
Challenge CLEF (Conference and Labs of the Evaluation Forum) 2025 
Modèle prédictif obtenu par apprentissage automatique (à implémenter)

Abstract The cyanobacterium Planktothrix rubescens Anagnostidis & Komarek (previously Oscillatoria rubescens DC ex Gomont) is present in several Italian lakes and it is known to produce cyanotoxins. The dynamics and toxin production of P. rubescens population in Lake Albano, a volcanic crater lake in Central Italy, has been studied for 5 years (January 2001-April 2005). Winter-



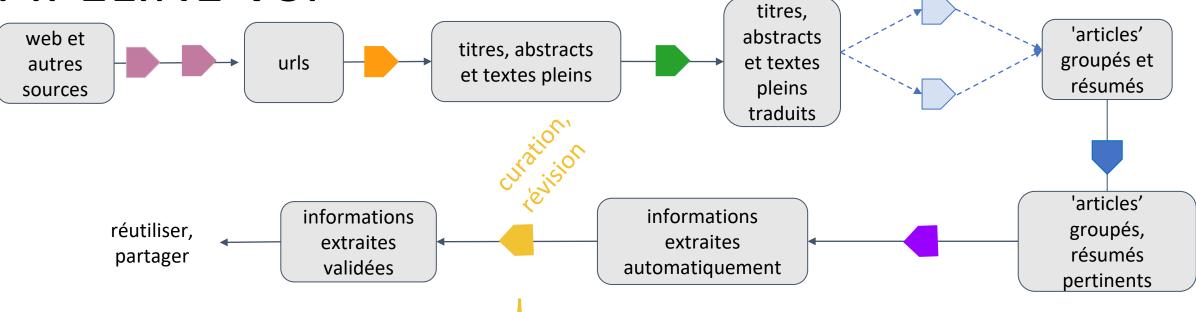
Abstract The cyanobacterium Planktothrix rubescens Anagnostidis & Komarek (previously Oscillatoria rubescens DC ex Gomont) is present in several Italian lakes and it is known to produce cyanotoxins. The dynamics and toxin production of P. rubescens population in Lake Albano a volcanic crater lake in Central Italy, has been studied for 5 years (January 2001-April 2005). Winter-

#### Article non-annoté

#### Modèle prédictif

### Article annoté automatiquement

(Execution modèle et enregistrement données dans BDD via API)



Comité éditorial via interface OVIDE – production manuelle des bulletins (scriptR)

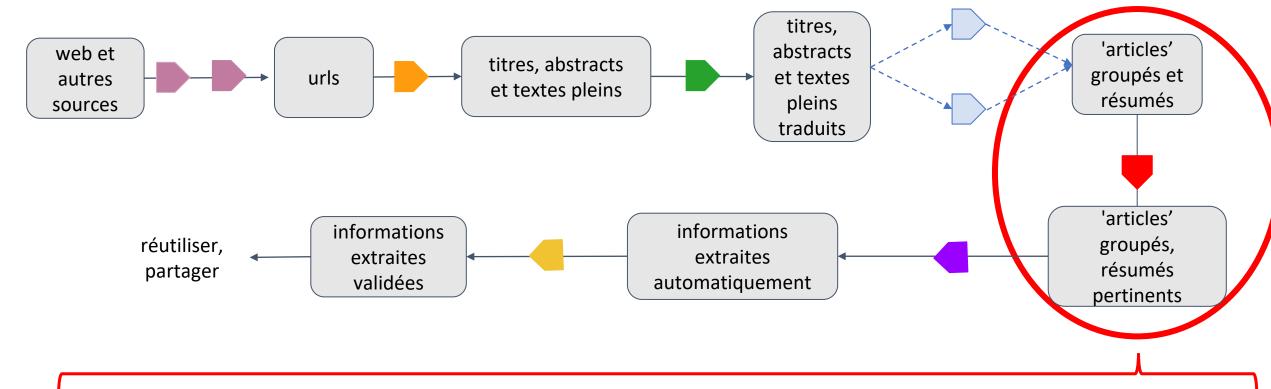
### **PERSPECTIVES**

- Production automatique des bulletins.
- Mise en place d'un moteur de recherche VSI directement associé à OVIDE production de bulletins à la carte Synthèse spatio-temporelle automatique (cartographie) automatique

•

Vectopole Sud, Montpellier, 18-09-2025

## CLASSIFIEUR Présentation

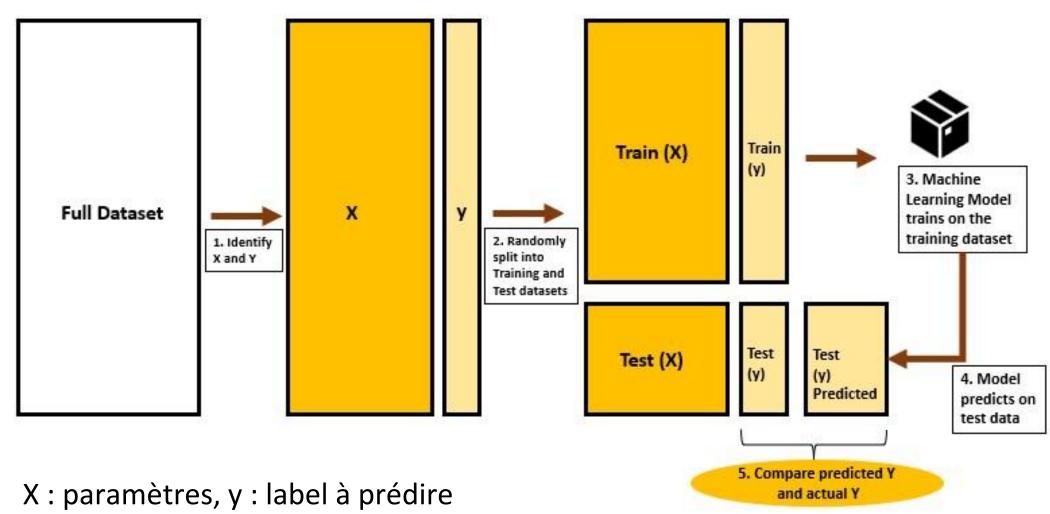


### Classification : sélection des articles pertinents

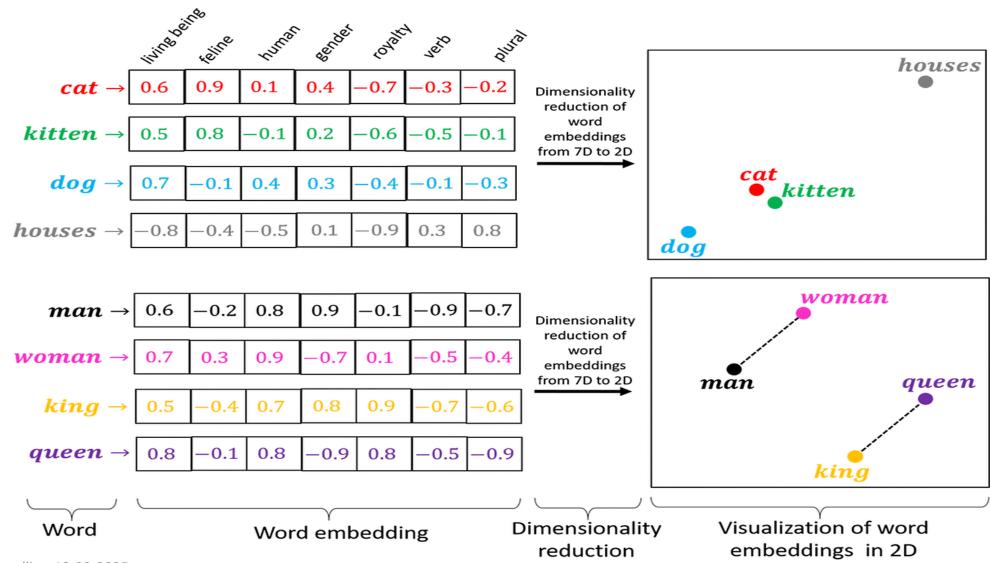
Environ 1000 articles collectés par semaine, dont ~15% d'articles pertinents. ~1 jour de travail par semaine par veilleuse —> objectif automatisation partielle Grandes quantités de données labellisées par les veilleuses



### Présentation : approche machine learning



### Entraînement: word embeddings



### CLASSIFIEUR Présentation : exemples

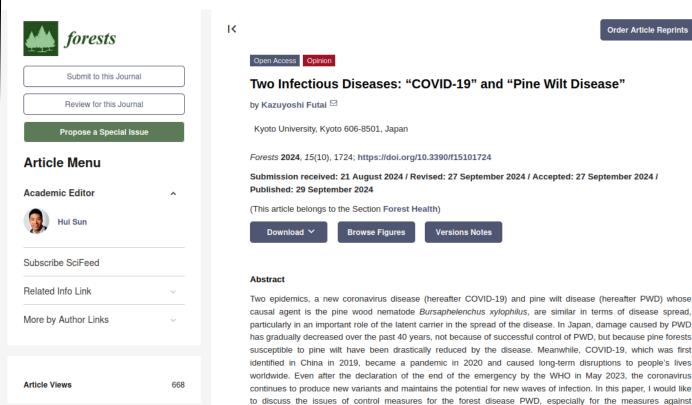


#### HLB: Se erradicaron 19 plantas positivas en dos localidades de la provincia de Corrientes

Lo hizo el Senasa en el marco de su Programa para la prevención de la enfermedad más destructiva de los cítricos.



**ID 217095 : Pertinent** 



Order Article Reprints Two Infectious Diseases: "COVID-19" and "Pine Wilt Disease" by Kazuyoshi Futai □ Kyoto University, Kyoto 606-8501, Japan Forests 2024, 15(10), 1724; https://doi.org/10.3390/f15101724 Submission received: 21 August 2024 / Revised: 27 September 2024 / Accepted: 27 September 2024 / Published: 29 September 2024 (This article belongs to the Section Forest Health) Two epidemics, a new coronavirus disease (hereafter COVID-19) and pine wilt disease (hereafter PWD) whose causal agent is the pine wood nematode Bursaphelenchus xylophilus, are similar in terms of disease spread. particularly in an important role of the latent carrier in the spread of the disease. In Japan, damage caused by PWD has gradually decreased over the past 40 years, not because of successful control of PWD, but because pine forests susceptible to pine wilt have been drastically reduced by the disease. Meanwhile, COVID-19, which was first identified in China in 2019, became a pandemic in 2020 and caused long-term disruptions to people's lives

asymptomatic carrier trees, with reference to the efforts implemented for COVID-19, which is more closely related to human society. This is because an asymptomatic carrier has been a blind spot in conventional PWD control

ID 214962 : Non pertinent



Vectopole Sud, Montpellier, 18-09-2025

### Entraînement : Pré-traitement des données

- Label 1 si pertinent, 0 sinon
- Concaténation du title + abstract + text (pas de lien / nom du site web)
- Tri des lignes vides —-> environ 16% des articles, pertinents par défaut
- « Accept cookies » ou « javascript cant download » gardés
- Si deux articles identiques/similaires ont des labels différents, ils sont relabellisés 1



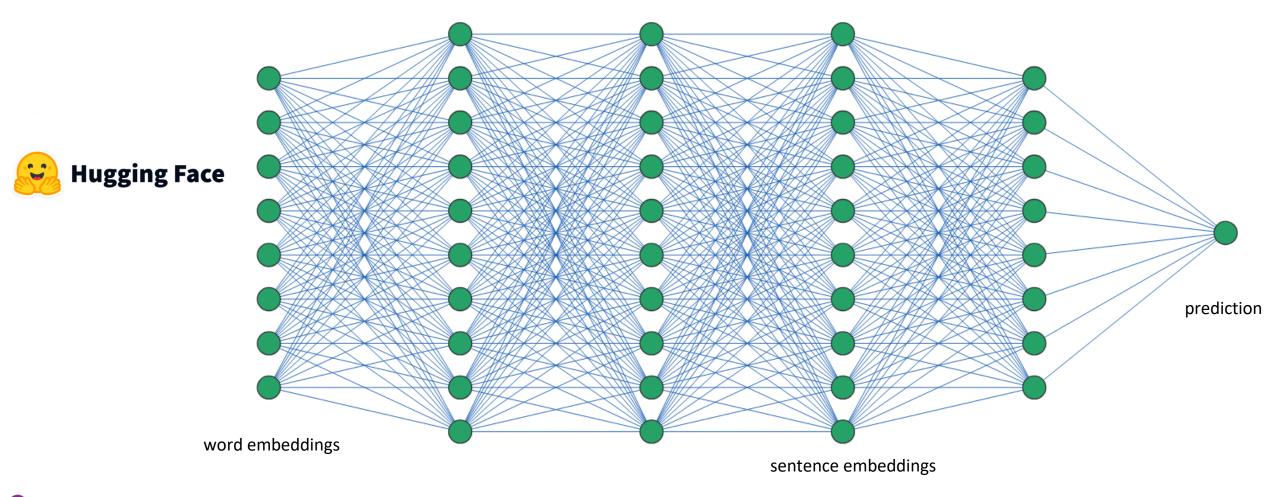
## CLASSIFIEUR Entraînement

1	id	date	lang	url	input	labels
2	213331	2024-09-16T05	es	https://www.elespanol.com/alica	Farmers in Alicante are furious over the government's "permissiveness" in the "avalanche of pests" from South Africa [SEP] The employers' association Asaja criticises Madrid's "passivity" in the face of phytosanitary problems affecting citrus fruits and calls for the Mercosur agreement to be reconsidered. Farmers in Alicante are outraged by the government's "permissiveness" in the face of the	0
3	213814	2024-09-16T05	es	https://www.todoalicante.es/ecor	Farmers demand cold treatment for citrus fruits arriving from Africa to stop the pests that ravage Alicante crops   TodoAlicante [SEP] La Unió estimates that the last six diseases reported in the last 15 years have led to a 40% increase in production costs. Sections Services We highlight You need to be registered to access this feature. Sharing options The Banski spider or the South African thrips are just two of the latest pests that are ravaging	1



### Entraînement : architecture testées

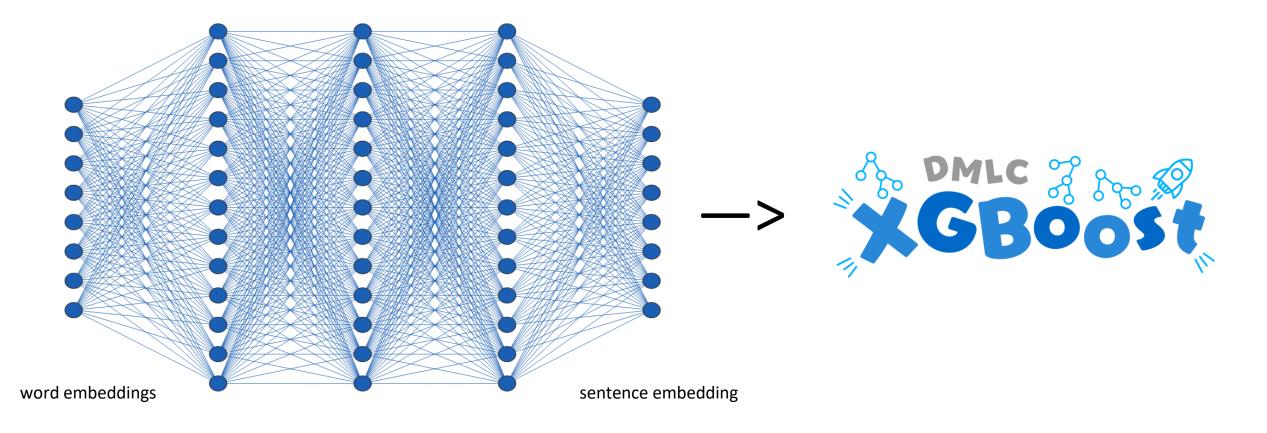
fine-tuning complet (100M paramètres)





### Entraînement : architecture testées

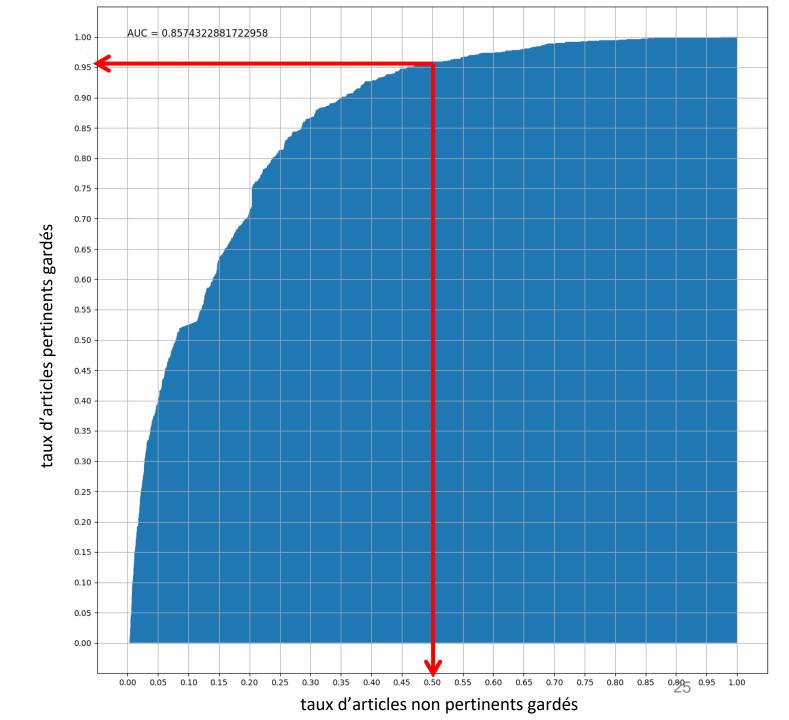
modèle de langue pré-entrainé (1B paramètres + entraînement local d'un modèle XGBoost)





Test: Résultats

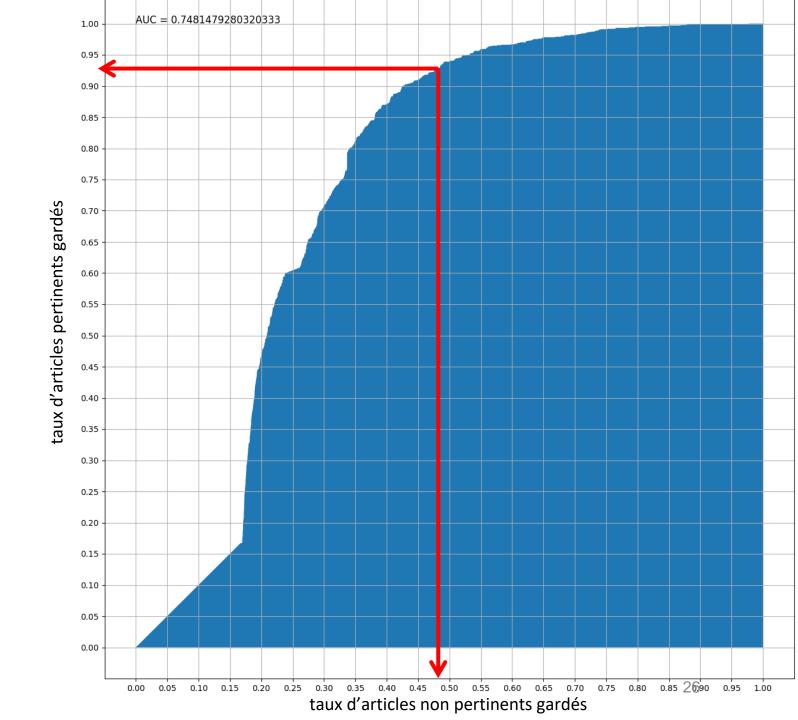
4% perte d'articles pertinents50% articles non pertinents éliminés



Test: Résultats

(lignes vides incluses)

6% perte d'articles pertinents —> 47% articles non pertinents éliminés



### Validation : Sélection du seuil discriminant

input =	label 🔽	prediction =	%_articles_pertinents_perdus =	%_articles_non_pertinents_éliminés =
NT+ Local Authorities & Construction [SI	1	0,0057	0,20%	0,50%
NT+ Local Authorities & Construction [SI	1	0,0057	0,41%	0,57%
Cantons - imTicker staht's [SEP] Top ne\	1	0,0069	0,61%	14,88%
Greenery Scanner, the electronic eye to	1	0,0076	0,82%	21,95%
Destination Canal du Midi   Facebook [S	1	0,0076	1,02%	22,17%
Salmonella in imported cucumbers. 18 p	1	0,0081	1,22%	23,49%
Multiscale Modelling of European Beech	1	0,0083	1,43%	24,18%
Transcriptomic response of citrus psyllid	1	0,0091	1,63%	28,29%
Webinar – Let's Talk About Planted Fore	1	0,0097	1,84%	29,29%
Emergency permit granted in the fight aç	1	0,0111	2,04%	30,39%
The State supports the agricultural trans	1	0,0126	2,24%	31,08%
Biological control of citrus pests: A syste	1	0,0149	2,45%	31,30%
Longhorn beetles from postharvest wood	1	0,0162	2,65%	,42%
Influence of Different Mango Phaeologic	1	0,0169	2,86%	36,86%
Flavescence dorée, regional tender wor	1	0,0187	3,06%	37,27%
Moderate Phosphorus Addition to Field-	1	0,0195	3,27%	37,58%
Undate of the delimited area for the pr				

Update of the delimited area for the pro Information Session: What is Xylella fa ToRREV virus in Poland - assessment

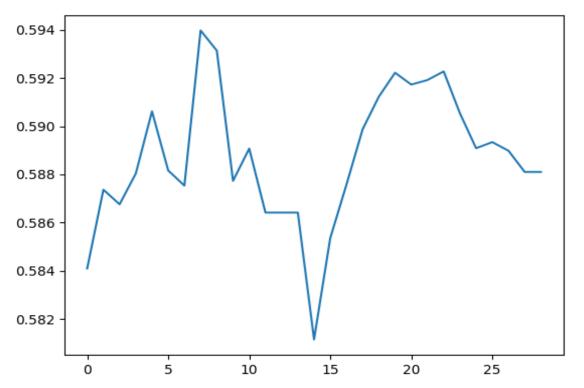
Acceptable de ne pas voir ces 13 articles pour réduire le volume à traiter d'un tiers ?



### Mise en production : amélioration continue

- Réorganisation de la BDD VSI, intégration à la pipeline VSI et à OVIDE
- Prévention d'une dérive du modèle (10% d'articles négatifs ré-échantillonnés)
- suivi des performances
- Ré-entrainement régulier

### Évolution de l'AUC en fonction du nombre de semaines





Mise en production : tracking des expériences

 Framework dédié au machine learning pour suivre les expériences (optimisation des hyper-paramètres, historique)



			Metrics	Parameters			
Run Name	Created =↓	Duration	AUC	embedding_model	learning_rate	n_estimators	run_date
	☑ 1 month ago	26.3s	0.83726205	NovaSearch/jasper_en_vision_language_v1	0.1	900	2025-07-11 17:51:38
charming-bass-54	✓ 1 month ago	29.7s	0.85459902	NovaSearch/jasper_en_vision_language_v1	0.1	900	2025-07-11 17:51:12
righteous-duck-449	☑ 1 month ago	35.3s	0.83560318	NovaSearch/jasper_en_vision_language_v1	0.05	900	2025-07-11 17:50:42
monumental-elk-366	✓ 1 month ago	38.0s	0.85618435	NovaSearch/jasper_en_vision_language_v1	0.05	900	2025-07-11 17:50:07
redolent-colt-661	✓ 1 month ago	22.3s	0.83645208	NovaSearch/jasper_en_vision_language_v1	0.1	600	2025-07-11 17:49:29
valuable-gnat-801	1 month ago	24.3s	0.85375877	NovaSearch/jasper_en_vision_language_v1	0.1	600	2025-07-11 17:49:06
sincere-goose-721	✓ 1 month ago	30.4s	0.83460830	NovaSearch/jasper_en_vision_language_v1	0.05	600	2025-07-11 17:48:42
legendary-snipe-945	1 month ago	31.8s	0.85457198	NovaSearch/jasper_en_vision_language_v1	0.05	600	2025-07-11 17:48:12

Vectopole Sud, Montpellier, 18-09-2025

### Mise en production : test grandeur nature juillet-août 2025

Classifieur comparaison veilleuses	Nombre d'articles
Vrai Positif	147
Vrai Négatif	369
Faux Positifs	424
Faux Négatifs	8

Semaines 28 à 35 (2025) ~ 999 articles /semaine



### Conclusion

- Données d'entraînement de qualité en quantité —-> bons résultats
- Meilleurs résultats avec des modèles pré-entraînés
- Cependant, modèles open source de taille raisonnable suffisants
- Résultats satisfaisants, passage en production définitif semaine prochaine ?

### Perspectives

- Ajout d'ONs pour voir le changement envisageable
- Périodicité de ré-entraînement
- Curation des données d'entraînement
- Amélioration du modèle open source
- Modèle plus gros ? Nouvelle architecture (RAG) ?
- Amélioration de la pipeline en amont



Isabelle Pieretti (animatrice du GT-VSI et veilleuse)



**Jean-Baptiste Louvet** (informaticien)



Marie Grosdidier (animatrice du GT-VSI et veilleuse)



**Simon Nicoux** (informaticien)



Sandy Dupérier (veilleuse)



Antoine Marullaz (Ingénieur en traitement automatique de données de télédétection et textuelles)





<u>isabelle.pieretti@cirad.fr</u> <u>antoine.marullaz@inrae.fr</u>