



# Intégration et développement d'approches de fouille de textes pour la veille épidémiologique en santé animale

### **Mathieu Roche**

Cirad, UMR TETIS, Montpellier

Montpellier, le 18 septembre 2025









### Context

### Text-Mining in Agriculture







(modestement, cela va sans dine) certaines activités, mais ils souhattent que la pleteforme continue et renforce ses activités au-delà de la fin de la petite convention CIRAD-FOFIFA qui a permis de faire émerger, et qui s'achèvere finalement à l'éta 2014. Ils sont apparemment prêts pour cela à s'investir dans la co-écriture d'un ne " pour le soumettre ensuité dans les nes pour le soumettre ensuité dans les les différents de la control de l

ne "pour le soumettre ensuire dans les bailleurs présents ou prévus sur le siche l'ocaux : Fonds régionaux de Développement mettent évantuellement en place : Alabotra : cella fait près de 3 ans qu'ils sont sgascar, ou tout autre bailleur bien disposé qui dans les prochains mois. Si un tel appul voyait ertinent que la recherche puisse accompagner fixtur monter au nouvrait loriture outre la fixtur monter au nouvrait loriture outre la

Corbesis Marc De Graaff zun Müch Hycenth Tinn Penelt Eric Baudenn Früderic Maudin Krishna Andries Medine. Christ Guillaumen Schuler Johannes Nyagumbo Isalisa Nisusinamhodal Leonord Trooré Karfin Mobal Hamisi Dilla. Adolwa Yuan Solomon. 2014. Understanding the impact and adoption of conservation agriculture in Africa: A multi-scale analysis. Agriculture. Cootsystem and Environment. 187: 115-170.

Quartile: Quatier, Suiet: AGRICULTURE, MULTIDISCIPLINARY / Quartile: Q1. Suiet: ENVIRONMENTAL SCIENCES / Quartile: Q1. Suiet: ECOLOGY

Understanding the impact and adoption of conservation agriculture in Africa: A multi-scale analysis

types (agriculture au sens large, pisciculture) auprès des membres de diverses OP, et particulièrement celles membres de la confédération

Nisumé: Conservation agriculture (CA) is increasingly promoted in more than two discusses of research and development investments. Through making past and on-pagin CA percenters in a set of case studies, this paper seeks to better understand the reasons for the limited adoption of CA and to assess where, when add for without CA through the control of the control of the control of the control of the destinguishes the following science of advantus. Held, farm, village and regions. CA has a potential to increase only within the following especially under conditions of errorit cantall and over the long-term as a result of a gradual increase of overall policity. The impact on farm increase with the practice of CA on some fields of the farm is far tess evolution. Also of participated in the Type of Earn is far tess evolution. We offer the control of the control of the control of the properties of the control of the description.

the non-adoption of CA. Smallholders have often short-term time

Industries | Wed Jul 23, 2014 9:54am EDT

Related: NON-CYCLICAL CONSUMER GOODS

#### Poland investigates suspected case of African swine fever in farm pigs

Polish local authorities said on Wednesday that preliminary tests have African swine fever (ASF) among farm pigs in eastern Poland near the

JOURNAL DU CAMEROUN.COM

BUSINESS SOCIÉTÉ SPORT CULTURE & LOISIRS

ite porcine fait des ravages dans l'extrême-nord

autorités de la région ont demandé aux populations de ne plus consommer la

Actualité Focus Conseils pratiques

The head of the Grodek co tests showed that ASF was

"We are marking the area," mats, were being taken.

Anna Wlodarczak-Semczul

Poland's chief veterinary of officer said a statement on

Visite du ministère de l'élevage et de la pêche dans le cercle de Tominian : La peste porcine africaine, une menace pour le développement

économique



ndu public le jeudi 24 juin 2010, le préfet du département du Diamaré dans la du Cameroun et dont la principale ville est Maroua, a demandé à la population nde de porc. Une décision qui fait suite à une épidémie sans précèdent de peste lans cette nartie du nave Le 02 iuin Nasseri Paul Réa avait officielle

[Drury and Roche, CEA'2019]





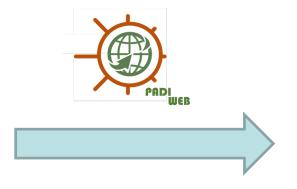




# PADI-web pipeline







### **Text-Mining**

 $\triangle$ 

Disease



Host



Location



Date



Number of cases



Symptom

AFRICAN SWINE FEVER

ANTIMICROBIAL RESISTANCE

AVIAN INFLUENZA

BLUETONGUE

BRUCELLA

CANINE DISEASE

CHIKUNGUNYA

CHOLERA

CLASSICAL SWINE FEVER

CRIMEAN-CONGO HAEMORRHAGIC FEVER

DENGUE

DISEASES OF BEES

FELIN DISEASE

FOOT AND MOUTH DISEASE

LEGIONELLA

LEPTOSPIROSIS

NIPAH

ONE HEALTH

PESTE DES PETITS RUMINANTS

RIFT VALLEY FEVER
SCHMALLENBERG VIRUS

SMALL RUMINANTS

TICK-BORNE ENCEPHALITIS

TROPILAELAPS

TULARAEMIA

UNSPECIFIED DISEASE

USUTU

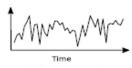
ZIKA

WEST NILE VIRUS

LUMPY SKIN DISEASE







Temporal Dimension



Dimension

.

Source Dimension





# PADI-web input



### Language and country labels



### Chile: Mysterious mass extinction of guano cormorants

ID: 8QTA0BGS1Q

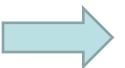
☐ Jun 4, 2023 · 
☐ Jun 5, 2023 · 
☐ Visit page

Source: latina-press.com

### Chile: Mysteriöses Massensterben von Guano-Kormoranen







B see less

Die Behörden untersuchen nun, was den Tod der Vögel verursacht haben könnte, ein Phänomen, das die Anwohner verängstigt hat (Foto: Agência Ambiental Pick-upau)

Datum: 04. Juni 2023 Uhrzeit: 15:49 Uhr Ressorts: Chile, Natur & Umwelt Leserecho: 0 Kommentare Autor: Redaktion

Im südamerikanischen Land Chile untersuchen die Umweltbehörden den Tod von Tausenden von Seevögeln und anderen Meeresbewohnern, deren Kadaver in der Region Coquimbo aufgefunden wurden. Nach Schätzungen der Landwirtschafts- und Viehzuchtbehörde (SAG) sind seit dem 26.





written permission of III IAP. Articles and reader reports identified by name do not necessarily reflect





# PADI-web pipeline





[Valentin et al., One Health 2021]



#### **PADI-web** Animal Health

#### **PADI-web** Plant Health

#### **PADI-web** Public Health







#### **Step 1: Data collection**

• Multlingual articles

#### Step 2: Processing

- · Text cleaning
- Language detection
- Translation into English

#### **Step 3: Document classification**

- Step 3.1: Relevant/irrelevant documents
- Step 3.2: Topic of documents

#### **Step 4: Sentence classification**

• Topic of sentences

### **Step 5: Information extraction**

- Disease
- Location
- Host
- Symptom
- etc.

### Step 6: Output for end-users

- Visualisations
- Notifications





# PADI-web process

■ Show source language



# Chile: Mysterious mass extinction of guano cormorants ID: 8QTAOBGSIQ





Chile: Mysterious mass extinction of guano cormorants

Chile: Mysterious mass extinction of guano cormorants ID: 8QTA0BGS1Q ☐ Jun 4, 2023 · 
☐ Jun 5, 2023 · 
☐ Visit page Source: latina-press.com KEYWORDS **✓ USER KEYWORDS** Keywords User dictionaries 🛱 disease (epidemiological 🗯 host information) **A** symptom various **⑥** location REPUBLIC OF CHILE REPUBLIC OF PERU SOUTH AMERICA AUTOMATIC KEYWORDS Named Entity Recognition AGÊNCIA LATINAPRESS NEWS & MEDIA COQUIMBO HUMBOLDT IAP ISLA LOBOS DE TIERRA SAG Extraction Institution Extraction People CHILEAN COMORIANS SOUTH AMERICAN aroup Extraction Location









**Jata annotation** 

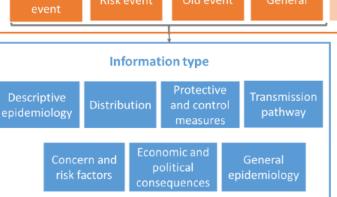
### Fine-grained classification

Current



Thematic Dimension





Corpus of sentences



- · Resolving of disagreements from dataset 3 (annotator A) · Re-annotation of dataset 1 and 2 (annotator A) · Aggregation of all datasets Final corpus
- Guidelines finalisation (1 hour) Step 2: Guidelines improvement (2 hours)

Step 1: Guidelines initialisation

(1.5 hours)

(3 hours)

- · Annotation of a third sample (dataset 3): 114 sentences (annotators A, B and C) · Agreement measurement
- · Second guidelines proposal (annotators A, B and C)

Third guidelines proposal (annotators A. B and C)

- Annotation of a 2<sup>nd</sup> sample (dataset 2): 209 sentences (annotators A and B), 194 sentences (annotators C and D)
- · Evaluation (calculation of metrics, discussion about the disagreements)
- · First guidelines proposal (annotator A)
- · Annotation of a first sample (dataset 1): 132 sentences (annotators A, B and C)
- · Evaluation (calculation of metrics, discussion about the disagreements)

		Inter-annotator agreement					
		Total agreements	Partial agreements	Disagreements	κ		
Step 1	Event type	29%	48%	23%	0.30		
	Information type	49%	43%	8%	0.53		
Step 3	Event type	87%	19%	4%	0.71		
	Information type	75%	22%	3%	0.78		





# PADI-web process







Alignments

X Semantic resource

✓ Automatic extraction

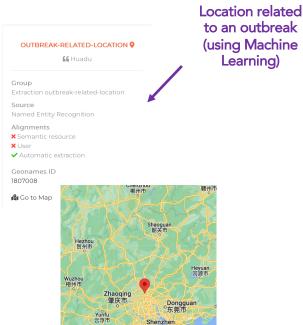
In any scase, he said, the veterinary board will continue to follow up on the plague until all

He pointed out that the ●outbreak of ● # African swine ● fever in ● Penang in ## #January

this year affected many 🌑 pig farmers in the state, and 🎈 🕈 Huadu village, as 🗠 one of the main 🜑

At the same time, the state government also arranged a second \$\infty\$ virus test after talking to \$\infty\$ \$\infty\$ pig farmers earlier, and found that many \$\infty\$ pig farms have no confirmed \$\infty\$ cases.

standard operating procedures are met before ----- is declared \$\ ASF-free.

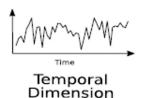






# PADI-web process



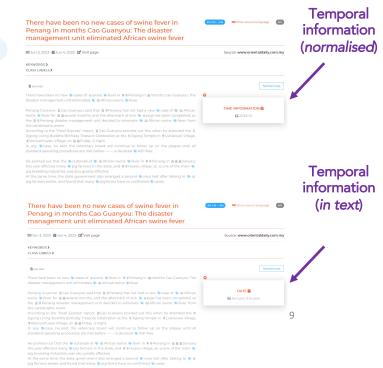


東方ONLINE

推荐 国内 国际 视频 财经 社会 地方▼ U玩食 娱乐

+更多

There have been no new cases of swine fever in **□□** Show source language Penang in months Cao Guanyou: The disaster management unit eliminated African swine fever ID: L7QCUXB2FE ■ Jun 3, 2023 - Jun 4, 2023 - Visit page Source: www.orientaldaily.com.my KEYWORDS > CLASS LABELS > Penang Governor 🎚 Cao Guanyou said that 🖟 🕈 Penang has not had a new 🗞 case of 🗞 🛎 African swine 🔷 fever for 🗎 🖺 several months, and the aftermath of sick 💜 🛣 pigs has been completed, so LOCATION • the 🛮 🕈 Penang disaster management unit decided to eliminate 🐎 🛎 African swine 🐃 fever from 66 Huadu village According to the "Pearl Express" report, III Cao Guanyou pointed out this when he attended the 9 Jigong Living Buddha Birthday Treasure Celebration at the ♥Jigong Temple in ♥Liutiaowei Village, In any scase, he said, the veterinary board will continue to follow up on the plague until all standard operating procedures are met before ----- is declared \$\ ASF-free. Alianments He pointed out that the Noutbreak of No African swine No fever in No Penang in A manuary this year affected many 🏶 pig farmers in the state, and 🎈 🕈 Huadu village, as 🗠 one of the main 🖠 X Semantic resource pig breeding industries, was also greatly affected. X User At the same time, the state government also arranged a second 🦠 virus test after talking to 💨 🐇 ✓ Automatic extraction pig farmers earlier, and found that many \$\infty\$ pig farms have no confirmed \$\infty\$ cases.







# PADI-web output



### Outputs

### Email **Notifications**



### Website consultation





### 14 new texts have been collected between 04 November, 2024 (00:00) and 05 November, 2024 > Go to search page

If this email is not fully working in your mail client, please click here to read it in PADI-web

#### CONTENTS

- Bluetongue (1) Africa (1)
- · African Swine Fever (6)
- · Europe (5) · Unknown (1)
- · Avian Influenza (5)
- Asia (1)
- Europe (3)
  North America (1)
- · West Nile Virus (2)
- Europe (1)
  Unknown (1)
- · Classification Information
- Keywords and Extracted Information

#### Bluetongue (1) ↑

#### Africa (1) ↑

#### Libyan Arab Jamahiriya (1) ↑



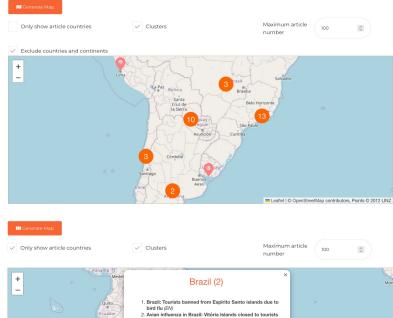
#### African Swine Fever (6) ↑

#### Europe (5) ↑

#### Italy (4) ↑



# Different view points to define an event





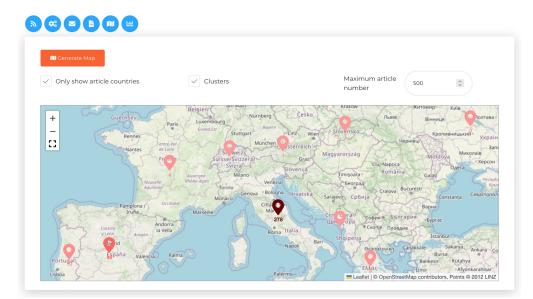


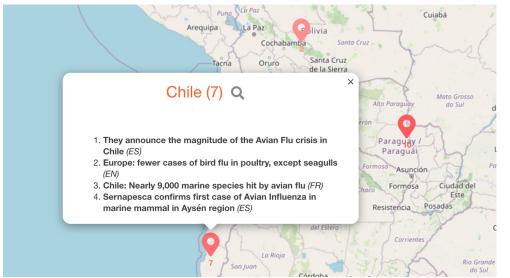




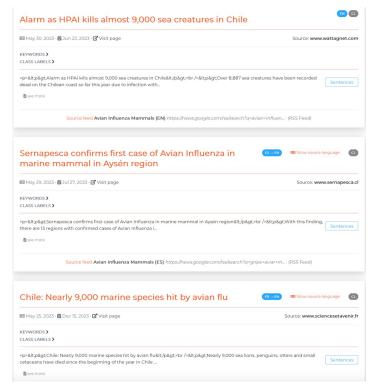




















[Trevennec et al., KES 2024]





**Carlène Trevennec** French Animal Health Epidemic Intelligence

### Provide a comprehensive view of emerging disease events

MUltisource Surveillance Tool (MUST) for Event-Based Surveillance (in the public health field)

#### **Public interface**

- > Extract disease event from various sources: official (IBS) and unofficial (EBS)
- > Space and time visualisation
- > Add complementary context information related to the outbreaks

#### **Advanced interface**

- > EBS optimization (validation of events, model training)
- > Epidemiological analyses







01/06/2025







#### Avian influenza:

Recent increase in reported infections in **mammals** and sporadic cases in humans Pandemic potential

Important for the French epidemic intelligence team – international monitoring of animal health, risk of introduction in France

### **MUST-Al integrates 3 different sources:**

- WAHIS: official source of public health notifications (database)
- **PADI-web**: unofficial sources that collect events from online media (Google news...)
- ProMED: expert networks in infectious diseases (mails and gradually a database)
- => collecting unstructured data from ProMED and PADI-web and extracting valuable information to make comparisons and complete official data









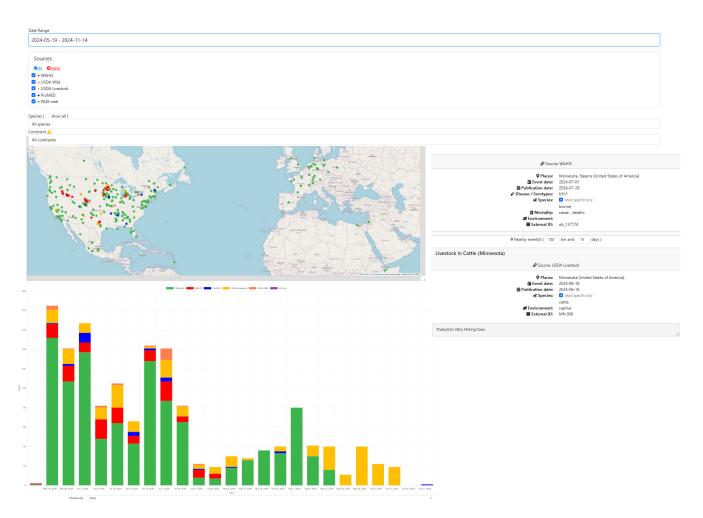


# MUST-AI INTERFACE

Open access tool:

https://must-surveillance.com/

















#### Strategy

- All
- SpaCy locations in Outbreak articles
- SpaCy locations in Outbreak articles and Current event sentences
- SpaCy locations in beginning of Outbreak articles
- PADI-web-specific locations
- SpaCy locations in beginning of articles



# Various types of extractions are proposed according on a number of criteria:

- Locations extracted at the beginning of the outbreak articles or not
- Locations extracted with SpaCy learned with labeled data or not

Relevant articles or not: an article is considered very relevant if its publication date is before the date of the report or official notification

[Trevennec et al., KES 2024]

#### **Fusion**

There are three main criteria to determine if two events are similar:

- **Geographical proximity:** if the distance is less than a reference distance
- Temporal proximity: if the difference in dates is less than a specified maximum of days
- **Species proximity:** if the Levenshtein distance of the species names is less than a specified maximum value



# Work in progress



 EpidGPT: A Combined Strategy to Discriminate Between Redundant and New Information for Epidemiological Surveillance Systems

Doc 1: African swine fever alert in Arad. Hundreds of pigs will be slaughtered								
The Local Center for Disease Control (CLCB) Arad approved a plan of measures								
to prevent the spread of the disease The Romanian authorities decided that								
Pig housing buildings and equipment will be disinfected								
Doc 2: An outbreak of African swine fever has been confirmed in a farm with								
over 300 pigs in Frumusani commune, Aluniş village, Arad county, county								
authorities announced today, Friday – August 18, which specified that all animals								
will be euthanized.								

Model	Prec%	$\mathrm{Rec}\%$	$F_1\%$
$\overline{BERT(PADIWeb_{novel}}$ +	81.49	79.63	80.55
BookCorpus)			
$RoBERTa(PADIWeb_{novel} +$	80.94	79.53	80.23
BookCorpus)			
$\overline{BioBERT(PADIWeb_{novel} + $	85.67	87.34	86.50
PubMedAbstracts)			
BioELECTRA(PubMed +	90.67	88.34	89.49
$PADIWeb_{novel})$			
Novel Epid GPT (WebCorpus +	92.89	91.23	92.05
$PADIWeb_{novel})$			

New information extraction with AI methods

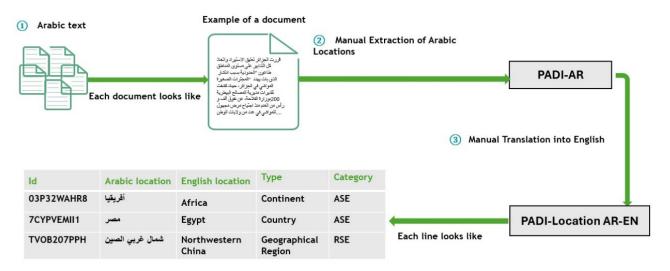
[Menya et al., NLDB'2024]



# Work in progress



Evaluation of the pipeline with Arab texts in Arabic language



Example of lines

Fatima Ezzahra El Houbri, Najlae Idrissi (Université Sultan Moulay Slimane, Faculté des Sciences et Techniques, Morocco) and Sandra Valentin (Cirad, TETIS)

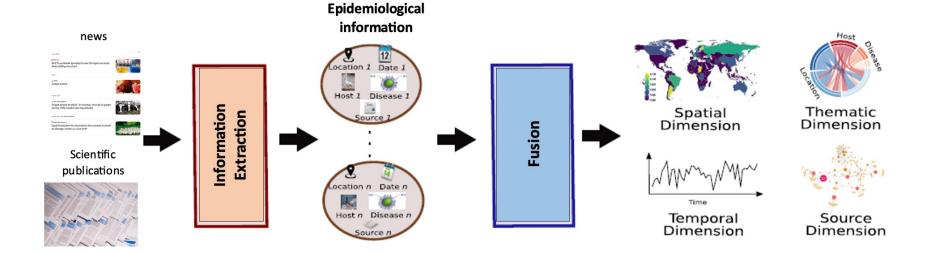
Integration of new NER system based on AI method (GliNER)

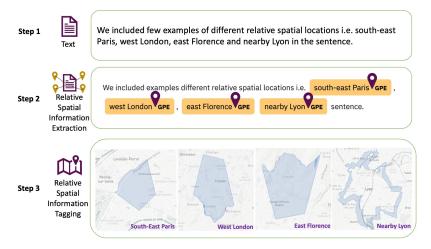




### Future work









### Future work







#### **Objectifs**

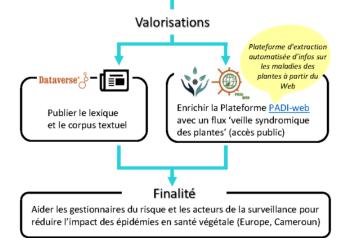
- · Détection précoce de nouvelles épidémies émergentes
- Détection précoce de nouvelles plantes hôtes

#### Méthodologie

Développer des approches de fouille de texte spécifiques à la veille syndromique végétale (cas d'étude européen et camerounais) mobilisant des compétences pluridisciplinaires :

- Santé : Epidémiologie végétale et animale
- Informatique : Fouille de texte, Intelligence artificielle





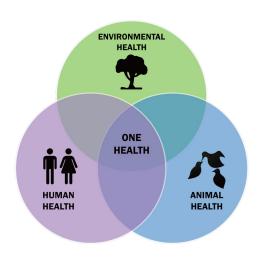
### In progress:

- Constitution of a corpus of syndromic surveillance. 2 types of queries: vocabulary of « mysterious » events & combination of « hosts » and « symptoms »
- Automatic augmentation of this dataset
- Proposition of a general One Health pipeline for syndromic surveillance



## **Future work**









Generecity

**Quality criteria** 

Other media





# PADI-web documentations







Animal Health: https://padi-web.cirad.fr/

Query 1

I would like to identify articles dealing with alert and preparedness for Avian Influenza in USA between 18/06/2023 and 01/09/2023

How many texts?

d. Other  $\square$ 

Objective:
Topics







18/11/2024

Objective:

Syndromic

surveillance /

**Topics** 

Animal Health: https://padi-web.cirad.fr/

#### Query 2

I would like to visualize countries classified as outbreak declaration of Avian Influenza in Europe between 01/12/2023 and 01/01/2024

#### How many texts?







Animal Health: https://padi-web.cirad.fr/

#### Query 4

I would like to explore (1) meta-data and (2) content of articles collected in March 2024 dealing with West Nile virus in Greece

#### What are other diseases mentioned?

- a. Zika 🗆 b. Covid-19 🗆
- c. Malaria ☐ d. Dengue ☐

#### What are the locations mentioned?

- a. Iraklio 🗆 b. Attica 🗆
- c. Tessaly 
  d. Samos





Objective: Information Retrieval

Animal Health: https://padi-web.cirad.fr/
How to retrieval texts with keywords?

I would like to identify texts dealing with mortality related to Avian Influenza?

Method 1: mortality as an indexed keyword

O Disease

O Published Benives

Notation

X

O Continent

M









I would like to identify the first sentence dealing with "economic and political consequences" in the article published the 7th of February 2024 and collected with the syndromic feed.

#### What is the first sentence dealing with economic and political consequences?

- a. Sentence 1 
  b. Sentence 4 
  c. Sentence 7







### Animal Health: https://padi-web.cirad.fr/ How to retrieval relevant period?

 $\begin{tabular}{ll} I would like to identify the month where Avian Influenza is the most frequent in 2024? \end{tabular}$ 











Animal Health: https://padi-web.cirad.fr/

#### Query 3

I would like to identify how many articles mentioned Avian Influenza and bovine in March 2024?

What are the original languages of the articles dealing with Avian Influenza collected on the 14th of March 2024 with poultry in the title?

#### How many texts?

- What is the language?
  - a. English 🗆
  - b. Spanish ☐ c. French ☐
- d. Chinese







#### How to implement notifications?

Other deficiency X and X

Step 2: You start the implementation of your notification







### PADI-web team







4 PhD students > 8 Msc students 1 post-doc

### Computer scientists



Mathieu Roche <sup>2</sup>
Senior Research Scientist



Sarah Valentin<sup>2</sup>
Research Scientist



Julien Rabatel IT development



# tetis TERRITIDIS ENVIRONMENT TÉL ÉDÉTETI

### **Epidemiologists**



Renaud Lancelot <sup>1</sup>
Senior Research
Scientist



Elena Arsevska <sup>1</sup> Research Scientist



Carlène Trevennec <sup>1</sup>
French Animal Health
Epidemic Intelligence



Isabelle Pieretti <sup>3</sup>
French Plant Health
Epidemic Intelligence







### References



### PADI-web v1 (2018)

Web monitoring of emerging animal infectious diseases integrated in the French Animal Health Epidemic Intelligence System. Elena Arsevska, Sarah Valentin, Julien Rabatel, Jocelyn De Goër De Herve, Sylvain Falala, Renaud Lancelot, Mathieu Roche. 2018. *PLoS ONE* 13(8): e0199960. https://doi.org/10.1371/journal.pone.0199960

### PADI-web v2 (2020)

PADI-web: A multilingual event-based surveillance system for monitoring animal infectious diseases. Sarah Valentin, Elena Arsevska, Sylvain Falala, Jocelyn de Goër De Herve, Renaud Lancelot, Alizé Mercier, Julien Rabatel, Mathieu Roche. 2020. *Computers and Electronics in Agriculture* 169: 105163. <a href="https://doi.org/10.1016/j.compag.2019.105163">https://doi.org/10.1016/j.compag.2019.105163</a>

### PADI-web v3 (2021)

PADI-web 3.0: A new framework for extracting and disseminating fine-grained information from the news for animal disease surveillance. Sarah Valentin, Elena Arsevska, Julien Rabatel, Sylvain Falala, Alizé Mercier, Renaud Lancelot, Mathieu Roche. 2021. *One Health* 100357. https://doi.org/10.1016/j.onehlt.2021.100357

### PADI-web Plant Health (2024)

**PADI-web for Plant Health Surveillance.** Mathieu Roche, Julien Rabatel, Carlène Trevennec, Isabelle Pieretti. 2024. In Proc. of CAiSE'24 - 36th International Conference on Advanced Information Systems Engineering, Springer, CAiSE Forum 2024: 148-156.



# PADI-web **Contact**



### Further questions and ideas on collaborations

padi-web@cirad.fr

https://www.padi-web-one-health.org











