



RÉPUBLIQUE  
FRANÇAISE

*Liberté  
Égalité  
Fraternité*



*maîtriser le risque |  
pour un développement durable*

# L'IA POUR RECONNAÎTRE DES EMPREINTES CHIMIQUES ?

09/06/2023

Ineris - 166139 - 2772091 0.1

# Sommaire

## 1. Position du problème

- a. Toujours plus de substances
- b. Toujours plus de données
- c. Pourquoi pas l'IA

## 2. Le POC « Empreintes environnementales »

- a. Les réalisations
- b. Le déroulé
- c. Des avancées réelles

## 3. Et maintenant

- a. L'industrialisation en question
- b. Quoi d'autre ?

# 1. Position du problème

# Des substances chimiques, encore des substances

## Quelques chiffres

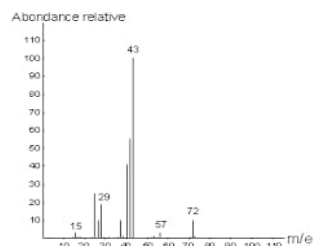
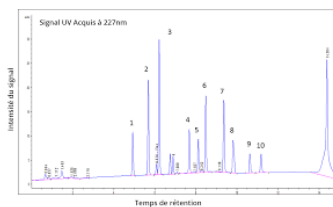
- 100 Millions de substances enregistrées dans le CAS\*
- 4000 de plus ... chaque jour !
- 30 000 à 50 000 substances pourraient se trouver dans les produits du quotidien



## Des substances et des méthodes

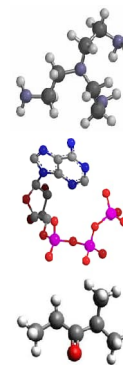
Pour caractériser des substances, il faut

- Des prélèvements d'échantillons,
- Des méthodes d'analyses,
- Des exploitations de résultats



## Tout ça prend beaucoup de temps et parfois ...

- Il y a urgence (Lubrizonl, Notre Dame de Paris...)
- On aimerait prévenir le risque plutôt que subir les conséquences



\*Chemical Abstract Service

# Des données, encore des données

## Sources des données

- Données de mesures en chromatographie et spectrométrie de masse (Ineris, Massbank, ACSM...)
- Données de constructeurs de matériels métrologiques
- Données de contexte environnemental
- Données de la littérature (rapports Ineris, US EPA...)

## Tout n'est pas sous Excel...

- Des formats complexes et difficiles à traiter par l'humain (MzML)
- Des données nichées dans des rapports et même dans des images
- Des données dans des formats propriétaires ou peu standards

```
1 22-04-20-Mix_15_Med_autoMMS_1_neg-1.mzML
2 <?xml version="1.0" encoding="utf-8" ?>
3 <indexedmzml xmlns="http://psi.hupo.org/mzml" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:schemaLocation="http://psi.hupo.org/mzml
4 http://psi.dev.info/files/mzml/xsd/mzml_1.0.xsd"
5 <unit xmlns="http://psi.hupo.org/mzml" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:schemaLocation="http://psi.hupo.org/mzml http://psi.dev.info/files/mzml/xsd/mzml_1.0.xsd" id
6 =22-04-20-Mix_15_Med_autoMMS_1_neg-1" version="1.0">
7 <cvlist count="2">
8 <cv id="MS" fullName="Proteomics Standards Initiative Mass Spectrometry Ontology" version="4.1.30" URI="https://raw.githubusercontent.com/PSI-P3/psi-ms-01/master/psi-ms.obo"/>
9 <cv id="UO" fullName="Unit Ontology" version="09-04-2014" URI="https://raw.githubusercontent.com/bio-ontology-research-group/unit-ontology/master/unit.obo"/>
10 </cvlist>
11 <fileDescription>
12 <fileContent>
13 <cvParam cvRef="MS" accession="MS:1000579" name="MS1 spectrum" value=""/>
14 <cvParam cvRef="MS" accession="MS:1000580" name="MSn spectrum" value=""/>
15 <cvParam cvRef="MS" accession="MS:1000127" name="centroid spectrum" value=""/>
16 <cvParam cvRef="MS" accession="MS:1000235" name="total ion current chromatogram" value=""/>
17 </fileContent>
18 <sourceFileList count="22">
19 <sourceFile id="AcqMethod.xml" name="AcqMethod.xml" location="file:///C:/Users/lestrému/Desktop/Francois/Travail Ineris/5_Non target
20 screening/data/Btalon_Avr20/NB0/22-04-20-Mix_15_Med_autoMMS_1_neg-1.d/AcqData">
21 <cvParam cvRef="MS" accession="MS:1001508" name="Agilent MassHunter nativeID format" value=""/>
22 <cvParam cvRef="MS" accession="MS:1001509" name="Agilent MassHunter format" value=""/>
23 <cvParam cvRef="MS" accession="MS:1000549" name="SHA-1" value="98FE07424944326c620db19246d491839Ea8"/>
24 </sourceFile>
25 <sourceFile id="BinPump1.og" name="BinPump1.og" location="file:///C:/Users/lestrému/Desktop/Francois/Travail Ineris/5_Non target
26 screening/data/Btalon_Avr20/NB0/22-04-20-Mix_15_Med_autoMMS_1_neg-1.d/AcqData">
27 <cvParam cvRef="MS" accession="MS:1001508" name="Agilent MassHunter nativeID format" value=""/>
28 <cvParam cvRef="MS" accession="MS:1001509" name="Agilent MassHunter format" value=""/>
29 <cvParam cvRef="MS" accession="MS:1000549" name="SHA-1" value="88c74b93081444696d3a3953348ea0b347207"/>
30 </sourceFile>
31 <sourceFile id="Contents.xml" name="Contents.xml" location="file:///C:/Users/lestrému/Desktop/Francois/Travail Ineris/5_Non target
32 screening/data/Btalon_Avr20/NB0/22-04-20-Mix_15_Med_autoMMS_1_neg-1.d/AcqData">
33 <cvParam cvRef="MS" accession="MS:1001508" name="Agilent MassHunter nativeID format" value=""/>
34 <cvParam cvRef="MS" accession="MS:1001509" name="Agilent MassHunter format" value=""/>
35 <cvParam cvRef="MS" accession="MS:1000549" name="SHA-1" value="285c1d4e26d451926f907fe79180d0E3E39b"/>
36 </sourceFile>
37 <sourceFile id="DefaultMassCal.xml" name="DefaultMassCal.xml" location="file:///C:/Users/lestrému/Desktop/Francois/Travail Ineris/5_Non target
38 screening/data/Btalon_Avr20/NB0/22-04-20-Mix_15_Med_autoMMS_1_neg-1.d/AcqData">
39 <cvParam cvRef="MS" accession="MS:1001508" name="Agilent MassHunter nativeID format" value=""/>
40 <cvParam cvRef="MS" accession="MS:1001509" name="Agilent MassHunter format" value=""/>
41 <cvParam cvRef="MS" accession="MS:1000549" name="SHA-1" value="8ba21925a203f299d4665f5d8f78297093a"/>
42 </sourceFile>
43 <sourceFile id="MethodPostData.xml" name="MethodPostData.xml" location="file:///C:/Users/lestrému/Desktop/Francois/Travail Ineris/5_Non target
44 screening/data/Btalon_Avr20/NB0/22-04-20-Mix_15_Med_autoMMS_1_neg-1.d/AcqData">
45 <cvParam cvRef="MS" accession="MS:1001508" name="Agilent MassHunter nativeID format" value=""/>
46 <cvParam cvRef="MS" accession="MS:1001509" name="Agilent MassHunter format" value=""/>
47 <cvParam cvRef="MS" accession="MS:1000549" name="SHA-1" value="8ba21925a203f299d4665f5d8f78297093a"/>
48 </sourceFile>
49 </unit>
50 </indexedmzml>
```

Exemple de fichier au format mzML

# L'IA pour nous aider à caractériser les substances

## Les méthodes d'IA semblent intéressantes, ça ...

- Ressemble à un problème d'apprentissage supervisé
- Traite des volumes importants de données
- Nécessite des temps limités de développement d'outils de calcul (vs. des modélisations physiques)
- Commence à être à la mode

## Quelques faits à prendre en compte

- Les équipes expertes de la chimie de l'environnement ne sont pas datascientists
- Les équipes sont déjà bien chargées
- Avant d'investir dans l'IA, il faut tester

## Que recherchions-nous ?

- Automatiser le processus de caractérisation des substances chimiques ou de leurs sources de contamination
- Augmenter et objectiver nos capacités d'expertise

## Et au passage...

- Structurer & valoriser nos données
- Apprendre à exploiter l'IA

## 2. Le POC « Empreintes environnementales »

# Fonctionnalités de l'outil développé

## Besoin #1

Obtention d'un score sur la probabilité d'identification d'une substance dans un échantillon

- ▶ Onglet **base de données** : affichage des échantillons ayant été intégrés dans la base, avec la possibilité d'intégrer de nouveaux échantillons et de réaliser une analyse
- ▶ Onglet **analyses** : recensement des analyses ayant été lancées, avec les paramètres associés, ainsi que les couples (échantillons, analyses)
- ▶ Onglet **résultats** : visualisation des résultats obtenus avec un accès aux scores définis par l'algorithme
- ▶ Onglet **dashboard** : visualisation des résultats détaillés par couple (échantillon, molécule)

## Besoin #2

Obtention d'éléments d'analyse et de comparaison objectifs qui permettent d'aider dans l'analyse et l'identification des sources de pollution dans un échantillon

- ▶ Onglet **exploration** : exploration de la base de données INERIS à travers ses différents degrés de liberté (projets, distance à la source, matrices et type d'exposition)
- ▶ Onglet **analyse** : mise en relief des données de la base INERIS pour obtenir des éléments de comparaison entre les différents échantillons qui composent cette base

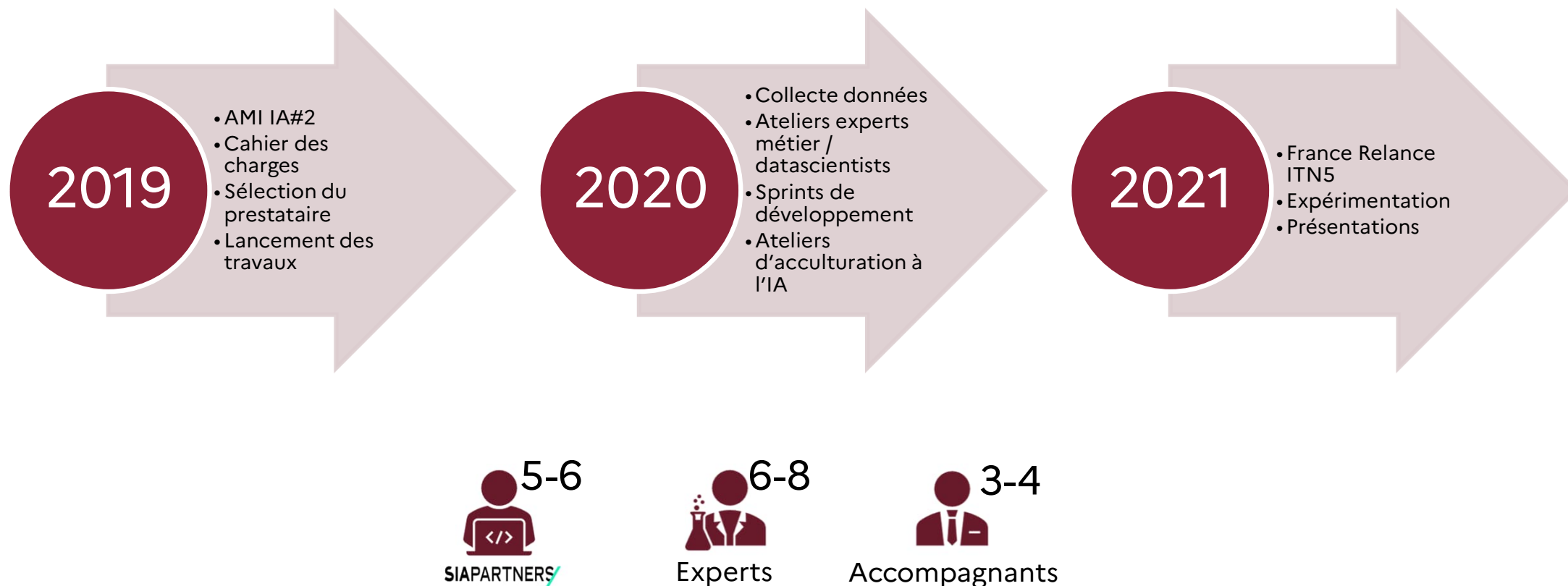
## Besoin #3

Détermination et visualisation en temps réel des sources de pollution à partir d'un échantillon d'air

- ▶ Onglet **dashboard** : visualisation de nombreux résultats graphiques et quantitatifs sur une période d'étude définie avec une modulation possible des graphes via les filtres
- ▶ Onglet **échantillon** : blacklistage des résultats aberrants et visualisation de l'état des tâches de chargement des données ACSM temps réel du SIRTA



# Organisation des développements



# Ce qu'on en retire

## Du point de vue métier

Des traitements de données plus fiables, plus exhaustifs et plus rapides

Un processus de caractérisation plus objectif : fixation de seuils, calculs de scores

Sensibilisation à la gestion des données en situation post-accidentelle notamment : structuration, gouvernance...

## Du point de vue applicatif

Une application rend les résultats partageables entre les acteurs

Il faut toujours bien garder à l'esprit que l'application reste un démonstrateur (moins frustrant)

Une application d'IA doit embarquer des fonctions type :

- apprentissage constant<sup>\*</sup>,
- traçabilité des modèles
- ...

## Du point de vue organisationnel

La mise en œuvre de l'IA permet :

- S'acculturer au métier pour le datascientist et réciproquement
- Catalyser l'ensemble du cycle de vie des données

<sup>\*</sup>Active learning

# 3. Et maintenant

# Une application à industrialiser (besoin #1)

## Développer une application dans les règles de l'art

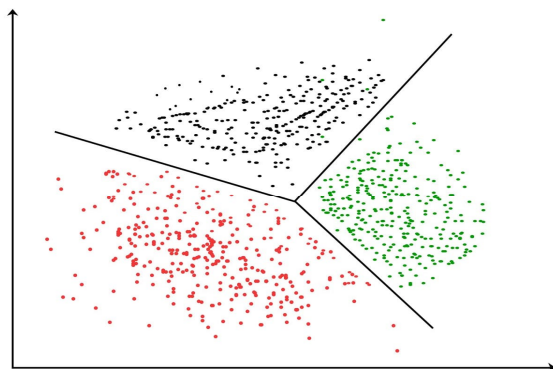
- Revoir le parcours utilisateur
- Traiter les questions de sécurité

## Maintenir l'application en conditions opérationnelles

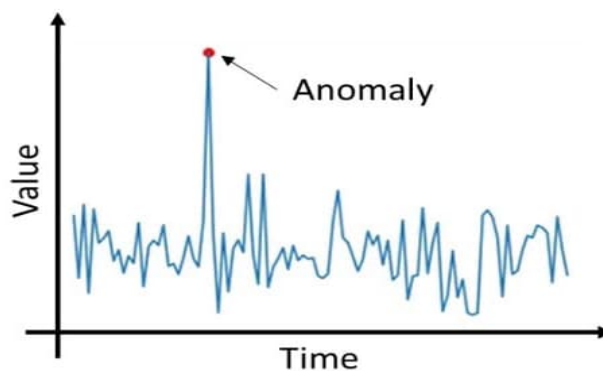
- Penser gouvernance des données : gestion des données d'apprentissage, active learning, traçabilité des résultats...
- Suivre les performances : temps de calcul, taux d'erreurs...
- Garantir son financement

# On pourrait imaginer aller plus loin

Clustering / groupes de substances



Recherches d'anomalies / suivi de sources



Utilisation de modèles prédictifs

PubChem