

# Common Data Elements.Artificial Intelligence CDE.AI

Nicolas Garcelon

► *Institut imagine – Plateforme data science*

Guillaume Assié, Anne-Sophie Jannot,  
Arnaud Sandrin, Xavier Tannier



**BNDMR**

Banque Nationale de Données  
Maladies Rares

► [bndmr.fr](http://bndmr.fr)

filère de santé  
maladies rares



**FIRENDO**

FILÈRE MALADIES RARES ENDOCRINIENNES



**Filnemus**

Filière Neuromusculaire

**fai2r**



**SORBONNE  
UNIVERSITÉ**



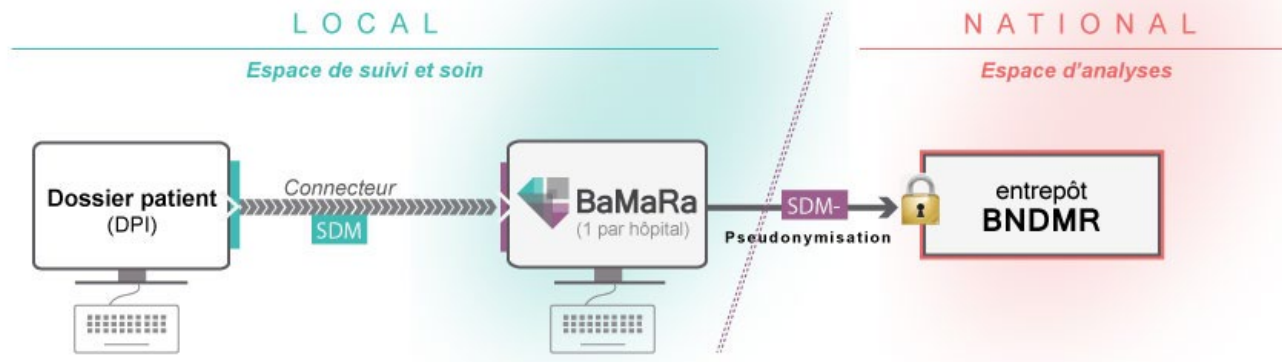
**institut  
imagine**  
GUÉRIR LES MALADIES GÉNÉTIQUES

**Université  
de Paris**



**ASSISTANCE  
PUBLIQUE**  **HÔPITAUX  
DE PARIS**

- ▶ Banque Nationale de Données Maladies Rares (BNDMR) : 60 items à saisir (dont ~20 obligatoires), 1,2M de patients



- ▶ Difficulté pour :
  - Exhaustivité des données à des fins de recherche (ex: 60% patients n'ont pas de signe clinique en HPO car donnée non obligatoire)
  - Qualité des données
  - Charge clinique : temps de saisie (création d'un dossier = 4 minutes par patient)
- ▶ Même problème pour toutes les bases recherches / registres

- ▶ Faciliter et accélérer la saisie des données pour la BNDMR
  - A partir du compte rendu d'un patient

## Compte-rendu

M. Z, born on February the 1st 2003

### Personal history:

- Left elbow fracture (2010)
- Low birth weight

### Familial history:

- Two sisters in good condition

### Treatments:

- Hydrocortisone 10: 1-0-0
- Fludrocortisone 50: 1-0-0

### Congenital adrenal hyperplasia:

- Diagnosed 5 days after birth by the systematic screening. sdfjisk
- Confirmed by mass spectrometry.
- Classical form with salt wasting
- No episode of adrenal crisis

## Remplissage automatique

Date of birth: XX/XX/XXXX

Diagnosis: YYYYYYYYYY

Familial history:  yes  no

Age at diagnosis: \_\_XX\_\_

## Correction/validation manuelle

Date of birth: XX/XX/XXXX

Diagnosis: YYYYYYYYYYYY




Familial history:  yes  no

Age at diagnosis: XX

- ▶ Reproductibilité sur une autre base :
  - Complications des glucocorticoïdes

# Méthode

*Rare diseases networks / expert centers*



**fai2r** **FIRENDO** **Filnemus**  
FILIÈRE MALADIES RARES ENDOCRINIENNES Filère Neuromusculaire

**>2000 unstructured medical records**

*Rare Diseases  
Common Data Elements  
(RD-CDE)*



**CDE.ai**

France Cohortes



**Inserm**

Secured data flow

*Pre-filled RD-CDE collection form*

**Human Validation**

**BNDMR**  
*Rare Diseases registry*



## 2 approches :

- ▶ Extraction par un modèle d'apprentissage profond
  - Annotation manuelle de comptes rendus pour « apprendre » la représentation d'une variable.
  - Application du modèle sur un nouveau texte pour identifier cette variable

Martine **PER** est née à **Paris LOC** le **17 Juin MISC** 1989 et habite maintenant à **Brest LOC**. Elle est atteinte d'un **syndrome auto-inflammatoire ORPHA** mais aussi de duplication partielle du bras long du chromosome 1. Après un voyage en **Alaska LOC**, il est atteint d'une **maladie due au virus Zika CIM10**. Une **incurvation du tibia LDDBfr** a été diagnostiquée

- ▶ Extraction à partir d'un pattern (expression régulière)
  - Définition d'un pattern pour extraire une donnée homogène :  
$$N[ée]e? *1e *([0-9][0-9][-/][0-9][0-9][-/][0-9][0-9][0-9][0-9])$$

# The Glucocorticoids Complication database

## Glucocorticoids Complications (Cushing syndrome)

- Endogenous
- Exogenous



Source: Wikipedia

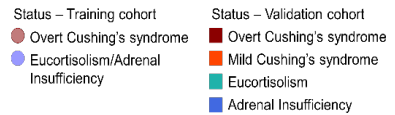
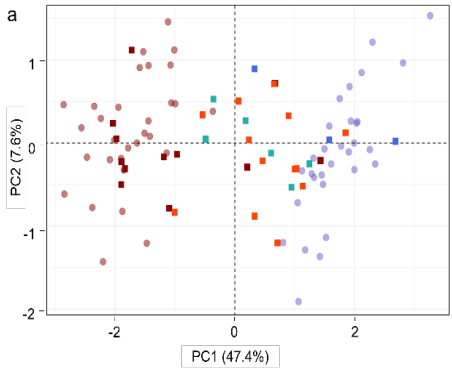
- Weight gain with faciotroncular adiposity
- Cutaneous and muscular atrophy
- Hypertension
- Diabetes mellitus
- Osteoporosis
- ....

## A need for molecular markers of Glucocorticoids Complications

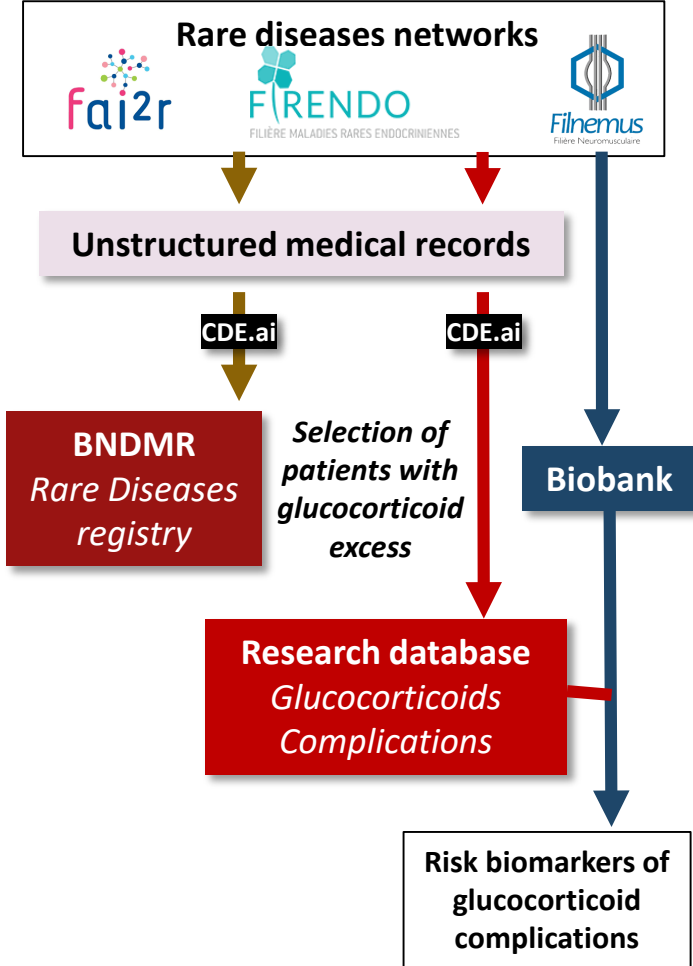
- Hormone assays not usable for exogenous Cushing
- Hormone assays do not predict individual complications
- Mild Cushing: should we treat or not?

### Markers:

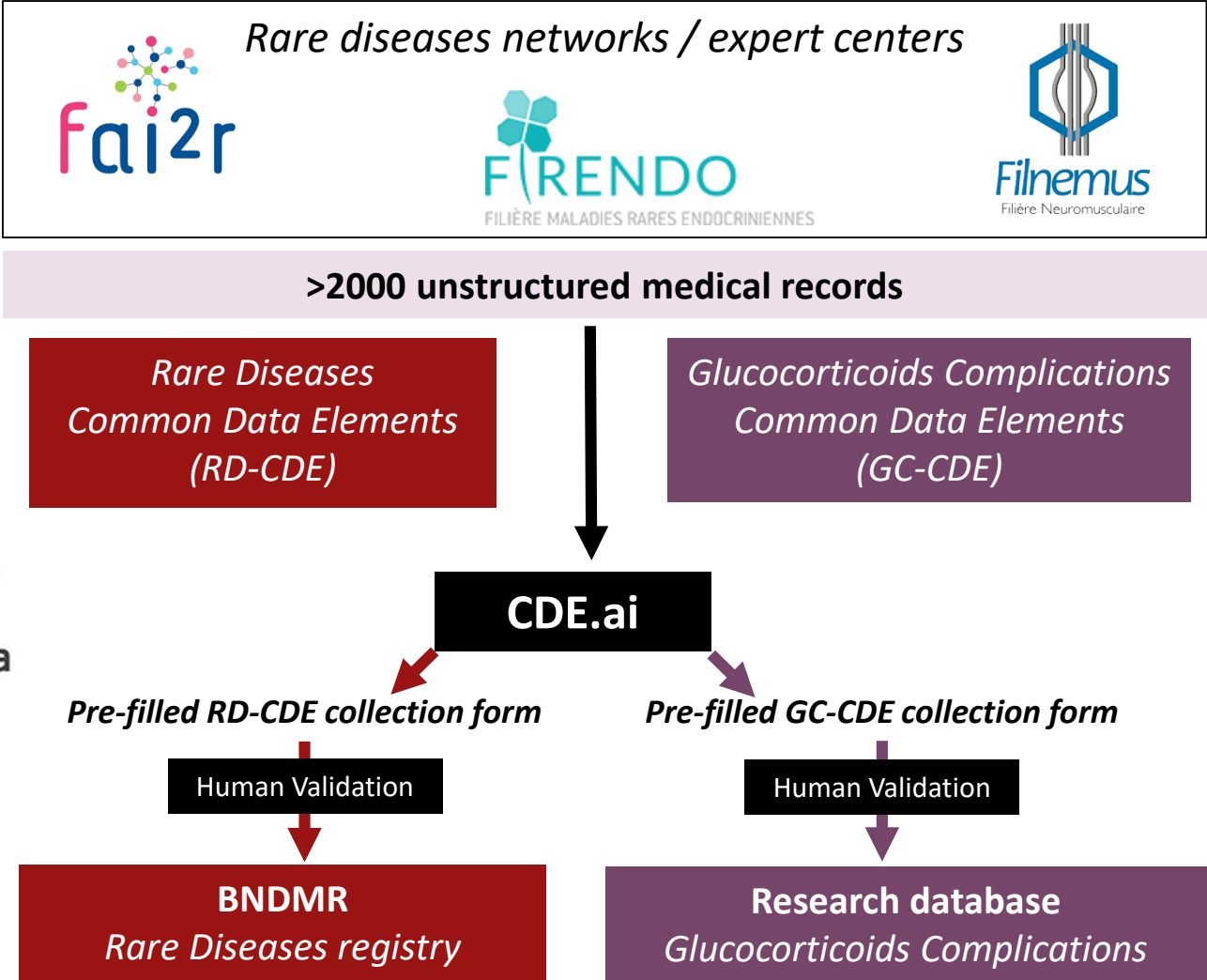
- GR variants
- Whole blood methylome



## GC-CDE Flowchart



# Méthode



France Cohortes



Secured data flow

*Annotation des comptes rendus*

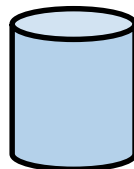
Martine **PER** est née à **Paris LOC** le **17 Juin MISC** 1989 et habite maintenant à **Brest LOC**. Elle est atteinte d'un **syndrome auto-inflammatoire ORPHA** mais aussi de duplication partielle du bras long du chromosome 1. Après un voyage en **Alaska LOC**, il est atteint d'une **maladie due au virus Zika CIM10**. Une **incurvation du tibia LDDBfr** a été diagnostiquée

*Remplissage du formulaire*

Date de naissance : 17/06/1989  
Diagnostic : syndrome auto inflammatoire  
Histoire familiale:  oui  non  
Age au diagnostic : 6 ans

*Collecte de la cohorte*

RD-CDE et GC-CDE enregistrements



Base de données

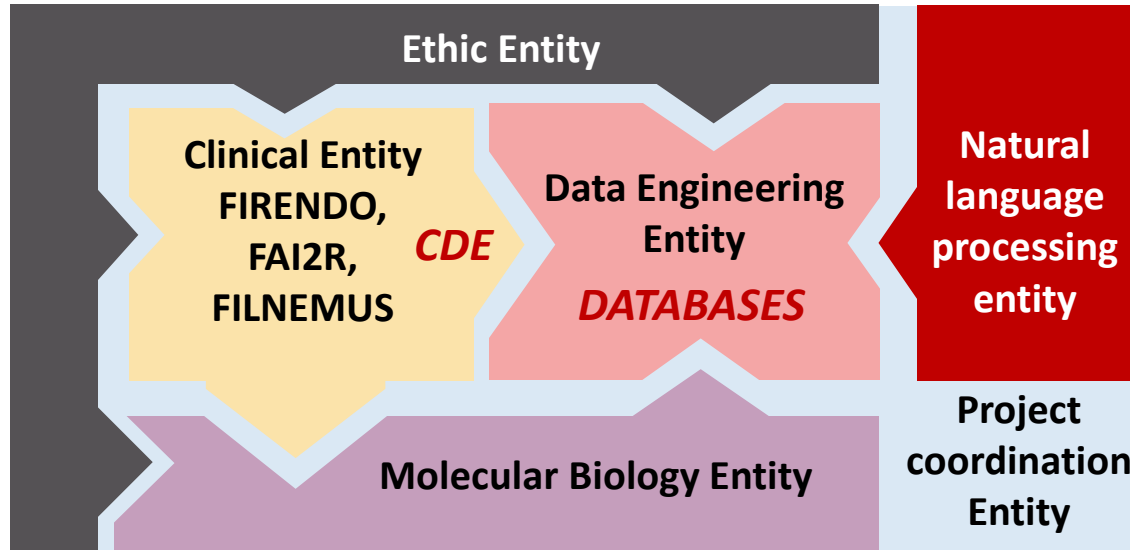
## 3 Niveaux d'évaluation

Précision de l'extraction des concepts

Qualité du remplissage des formulaires

Gain de temps





Un consortium  
multidisciplinaire

Team	Name	Affiliation
1	Assié G	Endocrine Dpt (Cohin APHP)
2	Jannot AS	Med Comput (HEGP APHP)
3	Sandrin A	BNDMR
4	Tannier X	LIMICS (Sorbonne U)
5	Garcelon N	Imagine, institut des maladies génétiques
6	Bertherat J	FIREENDO (APHP)
7	Belot A	FAI2R (HCL)
8	Attarian S	FILNEMUS (APHM)
9	Parfait B	CRB (Cochin, APHP)
10	Ragazzon B	Institut Cochin (INSERM)
11	Mamzer MF	Ethics lab ETREs (U Paris)

- ▶ **Conception et entraînement de CDE.AI**
  - Basé sur le **Rare Disease – Common Data Element**
  - Utilisation de la BNDMR, **infrastructure nationale** pour les maladies rares
  - Avec des spécialistes du **Traitement Automatique du Langage**
  
- ▶ **CDE.AI pour accélérer les bases recherche**
  - Sur les complications des glucocorticoïdes
  - Collecter des informations de biologie moléculaire avec GC-CDE
  - Réutilisable pour toutes les maladies rares et les questions de recherche.
  
- ▶ **L'importance de France Cohortes pour**
  - Héberger les dossiers médicaux pour l'entraînement de CDE.AI
  - Héberger la base de données des complications liées aux glucocorticoïdes
  - Favoriser l'avenir des collections de données sur les maladies rares avec CDE.AI
  
- ▶ **Impliquer les patients pour**
  - La réutilisation des données médicales hospitalières