# Data secondary use
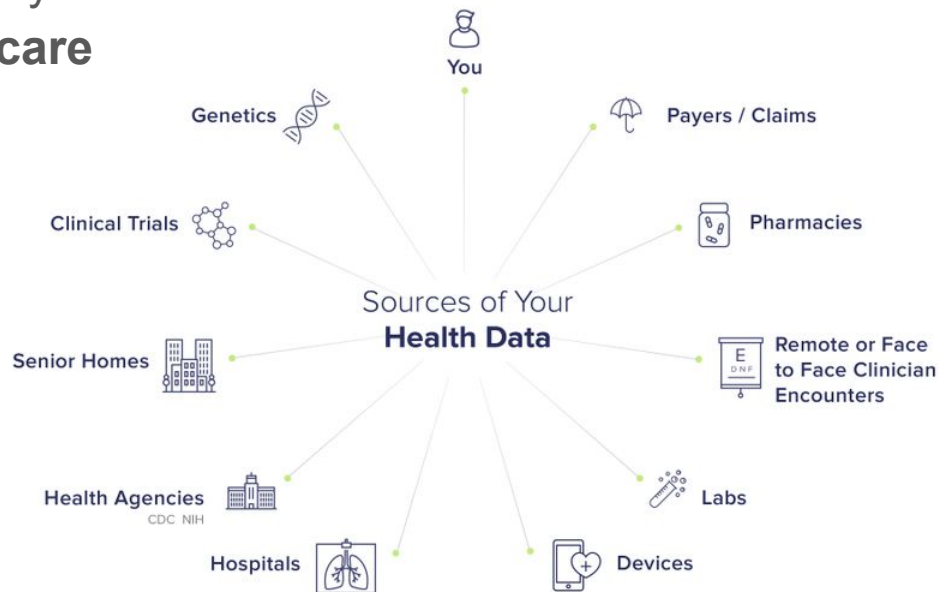
Master Public Health Data Science

Vianney Jouhet, MD, PhD
vianney.jouhet@u-bordeaux.fr

# Context

- More and more data are produced daily
- **Information technologies in healthcare**
  - **Reimbursement data**
  - **Electronic Health Record (EHR)**
  - **Biology**
  - **Radiology**
- **Research data**
  - **Clinical research**
  - **Epidemiology**
- Internet of things
- Omics

# Secondary use of biomedical data

1

Perspectives on Informatics

JAMIA

*White Paper* ■

## Toward a National Framework for the Secondary Use of Health Data: An American Medical Informatics Association White Paper
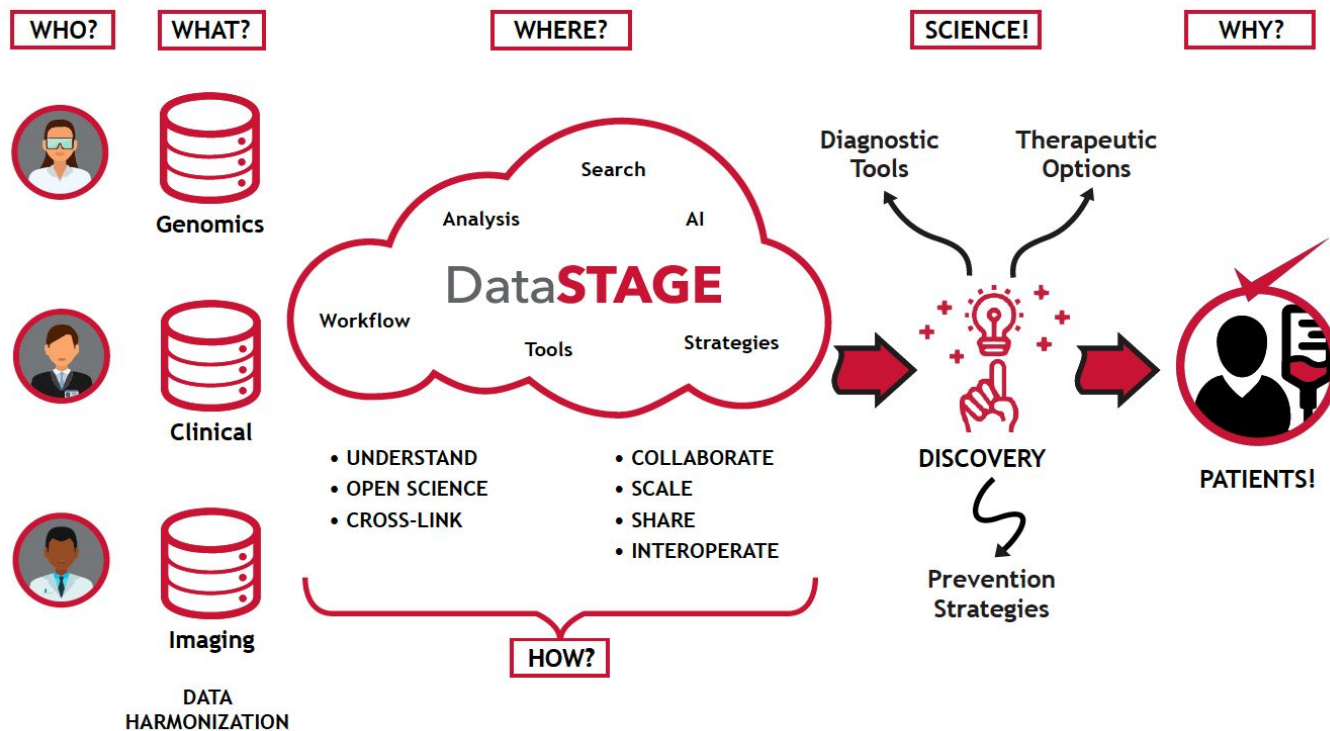
CHARLES SAFRAN, MD, MS, MERYL BLOOMROSEN, MBA, W. EDWARD HAMMOND, PHD, STEVEN LABKOFF, MD, SUZANNE MARKEL-FOX, PHD, PAUL C. TANG, MD, DON E. DETMER, MD, MA, WITH INPUT FROM THE EXPERT PANEL (SEE APPENDIX A)

# Secondary use of biomedical data

*"Secondary use of health data can enhance healthcare experiences for individuals, expand knowledge about disease and appropriate treatments, strengthen understanding about the effectiveness and efficiency of our healthcare systems, support public health and security goals, and aid businesses in meeting the needs of their customers"*

# Reusing research data DataSTAGE

# Database of Genotypes and Phenotypes (dbGaP)

- Archive and distribute NIH studies
    - Genotypes
    - Phenotypes
- Enable data retrieval
- Can not be computed directly
    - Encrypted files
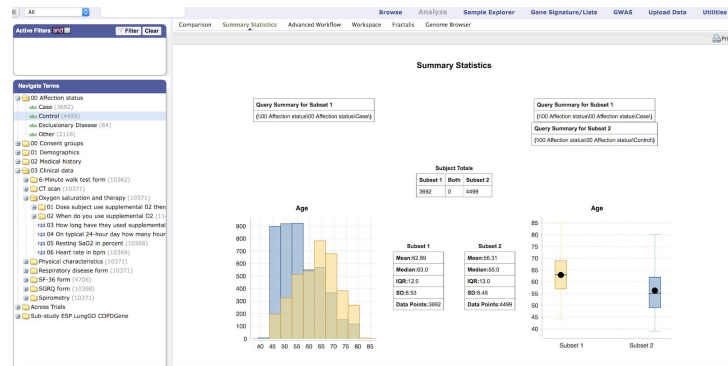    - XML dictionaries

# Architecture

# Clinical data Networks

# Clinical data Networks

# Clinical data Networks

# Secondary use and AI...



Mission VILLANI

DONNER UN SENS À L'INTELLIGENCE ARTIFICIELLE

POUR UNE STRATÉGIE NATIONALE ET EUROPÉENNE



HEALTH DATA HUB

MISSION DE PRÉFIGURATION

Une mission pilotée par Marc Cuggia (CHU Rennes), Dominique Polton (INDS), Gilles Wainrib (OWKIN) et rapportée par Stéphanie Combes (DREES)

# AI is a tool... not a goal

- User need comes before method choices
    - Identify domain expert
    - Define users needs
    - Choose methods
- Evaluate➜ There is no magic !

# Choose appropriate tools...

# Identify user needs...

Needs

Methods

# What can we do with these data ?

# Acute post-transfusion pulmonary edema

- Transfusion adverse event
- Should be declared by physicians
  - Underestimated
- Methods:
  - Simple free text search for concepts co-occurrence in a single sentence
  - Manual review of the sentence (fast human filtering)
  - Contextual validation of selected cases

# User Friendly UIs

CHU
Hôpitaux de **Bordeaux**

Hémovigilance

Identification

cossins

●●●●●●●●●●

Se connecter

Signaux détectés    Signaux validés

## Signaux d'OAP post-transfusionnel détectés :

- Signal n°11 détecté le 21/02/2018
- Signal n°12 détecté le 16/02/2018
- Signal n°13 détecté le 03/01/2018
- Signal n°14 détecté le 22/01/2018
- Signal n°15 détecté le 23/04/2018
- Signal n°17 détecté le 23/04/2018
- Signal n°18 détecté le 17/04/2018
- Signal n°24 détecté le 15/05/2018

*Patiente entree pour OAP post transfusion chez une patiente GIR 2 souffrant d Alzheimer et vivant en EHPAD*

Ignorer    Valider
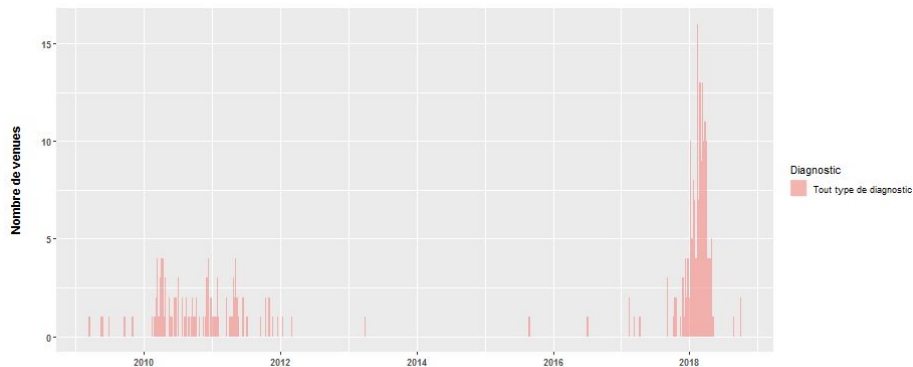
**DetectTACO**

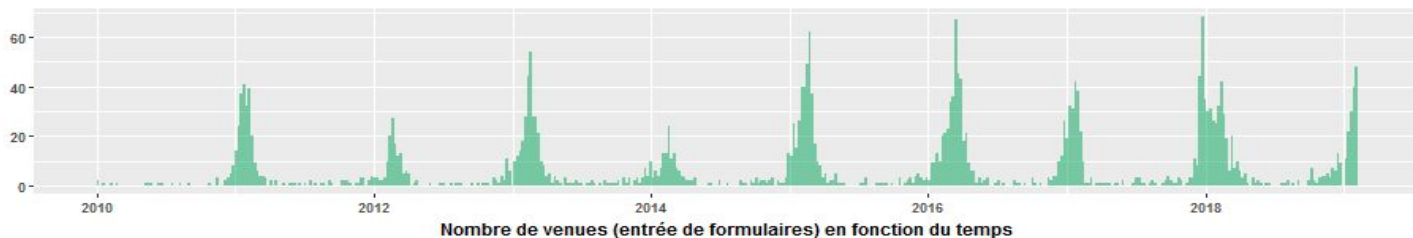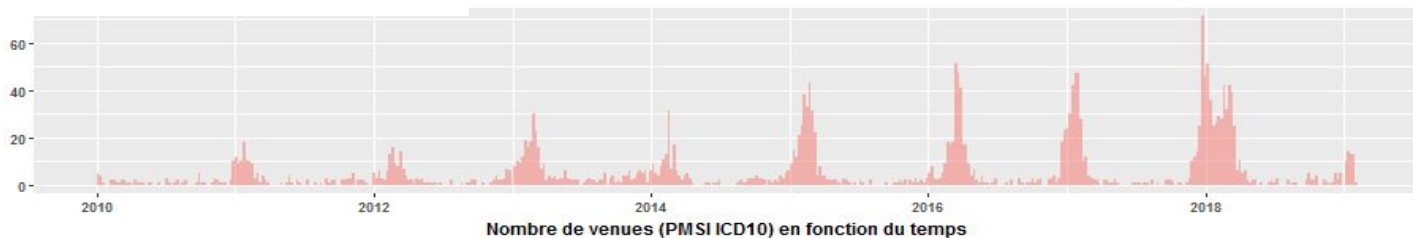Outil de détection des OAP post-transfusionnels.

# Results

- Without secondary use ⇒ 1 case identified (3 months)
- Using data
  - 102 « possible cases » detected
  - 8 cases validated
  - Cluster identification
  - Apply correction mesure

# Epidemic monitoring

# Phenotyping



Haendel, M.A., Chute, C.G., Robinson, P.N., 2018. Classification, Ontology, and Precision Medicine. N. Engl. J. Med. 379, 1452–1462. https://doi.org/10.1056/NEJMra1615014

# Phenotyping

Orphanet Journal of
Rare Diseases

**RESEARCH**　　　　　　　　　　　　　　　　　　　　**Open Access**

# Next generation phenotyping using narrative reports in a rare disease clinical data warehouse

Nicolas Garcelon[1,2,13*], Antoine Neuraz[2,3], Rémi Salomon[1,4], Nadia Bahi-Buisson[1,5], Jeanne Am
Capucine Picard[1,8,9], Nizar Mahlaoui[1,8,10,11], Vincent Benoit[1], Anita Burgun[2,3,12] and Bastien Ranc

Contents lists available at ScienceDirect

## Journal of Biomedical Informatics

journal homepage: www.elsevier.com/locate/yjbin

## A clinician friendly data warehouse oriented toward narrative reports: Dr. Warehouse

Nicolas Garcelon[a,b,*], Antoine Neuraz[b,c], Rémi Salomon[a,d], Hassan Faour[a], Vincent Benoit[a],
Arthur Delapalme[a], Arnold Munnich[a,e,f], Anita Burgun[b,c], Bastien Rance[b,g]

Elizabeth G. Phimister, Ph.D., *Editor*

# Classification, Ontology, and Precision Medicine

Melissa A. Haendel, Ph.D., Christopher G. Chute, M.D., Dr.P.H.,
and Peter N. Robinson, M.D.

# Need for semantic standardisation

- Identify information meaning uniquely
- Extract concepts from free text
- Organise information depending on its meaning
    - Medication
    - Phenotypes
    - Disease
- Drive data Visualization through external knowledge
- Based on Semantic resources
    - Terminologies
    - Ontologies

# Re-using hospital data

# Hospital Information systems

- Multiple distributed medical applications
    - Electronic Health record
    - Administrative data
    - Radiology
    - Biology
    - Etc…
- Build for healthcare purpose (not for secondary use)
- Reimbursement data production

**Lead to data silos**

# Hospital Information Systems

# Interoperability

- Standards (IHE, HL7, CDISC)
  - Messages (HPRIM, HL7)
  - Format
- Information exchange
- Depends on vendors implementation  +++
- Complex to implement

**Recently standard HL7 / FHIR**

https://www.hl7.org/fhir/

# Issues

Separated data (silos)

Data Heterogeneity

Bordeaux University Hospital

- 100+ Applications
- 10 000+ Tables

# Clinical data warehouse (CDW)

# ETL - Extract

- Extract data from Hospital Information System
    - Need knowledge of data model
    - This usually needs to be reverse engineered
- Another solution
    - Use exchange standard
    - Leverage standardised messages
    - Necessitate to listen data flows in real time
    - FHIR may offer easier solution for data extraction

# ETL - Transform

Transform data from HIS

- Technical transformation
    - Relational model to CDW information model
- Semantic transformation
    - Harmonize terminologies
    - Annotate information (meta modeling)
    - Leverage Ontologies and terminologies
    - May benefit from multi terminology server

# ETL - Load

Load data into the EDS information model

- Trade of
    - Update data (add, remove, update)
    - Drop and replace

# Example Bordeaux University Hospital CDW

# CDW: data integrated

# CDW: data integrated

**1 650 454**

Patients

**12 098 270**

Venues

**1 237 180 900**

Observations

Forms      Drugs      Lab tests      Discharge summaries      Radiology reports

# Data access and data privacy

# Driving Principles

- Security et traceability
  - Data are processed in situ (move  algorithm not data).
- Transparency
  - Patient information (General Data Protection Regulation)
  - Open source / Open data / open science

- Optimise data availability for secondary use and data privacy

# Confidentiality, integrity, availability

- Trade of between confidentiality, integrity, availability
- A system with high security will limit availability
  - Highest security ⇒ offline storage (not accessible at all !!!)
- A system with high availability may cause privacy breach and data leak
  - The most available ⇒ Open data available for all (no authorization, no authentication!!!)
- High integrity is time consuming
  - Perfect and complete data is unreachable (data will never be available)

# Data privacy risk levels

4 levels

1. Dictionaries (biologie, formulaire etc…)
2. Category counts
3. Query with patient numbers (obfuscation)
4. Detailed data
   a. Identifiers (HIPAA -  Health Insurance Portability and Accountability Act),
   b. Pseudonymised

# Goals of data processing

- Within structure scope
    - Health care quality, vigilances, evaluation ...
- Research, studies
    - Clinical research
    - Epidemiology
    - Research databases
    - Biobanks

# What can be done...

|  | Routine activity | Research, studies |
|---|---|---|
| Level 1 (dict) |  |  |
| Level 2 (cat counts) |  |  |
| Level 3 (counts) |  |  |
| Level 4 (detailed) |  |  |

# What can be done...

|  | Routine activity | Research, studies |
|---|---|---|
| Level 1 (dict) | Open data | Open data |
| Level 2 (cat counts) |  |  |
| Level 3 (counts) |  |  |
| Level 4 (detailed) |  |  |

# What can be done...

|  | Routine activity | Research, studies |
|---|---|---|
| Level 1 (dict) | Open data | Open data |
| Level 2 (cat counts) | Open data | Open data |
| Level 3 (counts) |  |  |
| Level 4 (detailed) |  |  |

Co-occurrence matrix (Natural language processing and structured data)
- Enables embeddings
- Train models

# What can be done...

|  | Routine activity | Research, studies |
|---|---|---|
| Level 1 (dict) | Open data | Open data |
| Level 2 (cat counts) | Open data | Open data |
| Level 3 (counts) | Authenticated user / widely available | Authenticated user / widely available |
| Level 4 (detailed) |  |  |

# What can be done...

| | Routine activity | Research, studies |
|---|---|---|
| Level 1 (dict) | Open data | Open data |
| Level 2 (cat counts) | Open data | Open data |
| Level 3 (counts) | Authenticated user / widely available | Authenticated user / widely available |
| Level 4 (detailed) | Authorized user / specific missions | Authorized user / Patient consent |