### Semantic Clinical Data Warehouse in Rouen University Hospital, Normandy, France





#### Prof. Stéfan Darmoni, MD, PhD Badisse Dahamna, Romain Lelong, Ivan Kergourlay, Kevin Billey, Julien Grosjean

Department of Biomedical Informatics (DBI), Rouen University Hospital, Normandy, France & LIMICS INSERM U1142, Sorbonne University, France



# RUH DBI

- Rouen University Hospital, Normandy
- N =15 (vs. 100 in Harvard Medical School)
  - Heterogeneity of competence ; double competence
  - Two medical informaticians
  - Five engineers
  - Three medical librarians (CISMeF, LiSSa)
  - Two/three medical residents in public health + one GP resident
  - Four PhD applicants

In 2019, RUH DBI became part of the LIMICS, main French lab in biomedical informatics

# Definition

- Health Data Warehouse (HDW) are computer tools allowing collection, integration, then data mining. These data are extracting from various clinical information systems: EHR, LIS, RIS/PACS, CPOE, etc...
- Agregation of a maximum of data available from patients
  - First step, HDW from DRG data
- No matter what is the Clinical Information System (CIS)
  - As far as possible, retrieving all the data from CIS, no matter old they are; going back to 2000 ≠ Paris Hospitals (APHP)
- ➔ Allowing to link information and to select as precisely as possible the right patients



# **HDW Objectives**

- To optimize DRGs by semiautomatic detection of atypical profiles between coding and HDW data (DBI)
- To Improve clinical research thanks to feasibility studies prior to clinical trials & optimization of inclusions
- Detection specific patients profiles
  - e.g. patients frequently admitted to the emergency department
- To create and maintain epidemiological cohorts and registries
- To detect specific adverse events
  - Vigilances, iatrogeny, cross infections
- To create and follow-up quality indicators
- To Develop and assess computer aided decision support systems (CDASS)

# État des lieux des EDS

Dans le monde :

• i2b2 🔰 (Harvard MS)

### En France :

- ehoc (CHU de Rennes)
- Dr Warehouse 🖳 (Necker, Foch Paris)
- Consorregional Continuum Soins Recherche



Lave





- 86 termino-ontologies included in 32 languages
- 2 M concepts in English; 1.2 M in French
  - $\approx$  145 K concepts in French in UMLS (2018AB) vs.  $\approx$  445 K in HeTOP (x2.98)
- Over 100 million RDF triplets (2014) big data +++

ECMT Semantic Annotator (NLP & Deep Learning)

### Layer 2

- Based on HeTOP; 50 chosen KOS out of 75 (no interface terminologies)
- 11.8 M health documents
- Processing time: 22-24 hours (two servers 1 To; one with 196 cores and the second with 144 cores
- 5.2 G medical concepts extracted; 2.6 G after filtering

Siefridt C, Grosjean J, Lefebvre T, Rollin L, Darmoni S, Schuers M. Evaluation of automatic annotation by a multi-terminological concepts extractor within a corpus of data from family medicine consultations. Int J Med Inform. 2020 Jan;133:104009. doi: 10.1016/j.ijmedinf.2019.104009. Epub 2019 Nov 1.

Neveol & coll. Clinical Natural Language Processing in languages other than English: opportunities and challenges. Journal of Biomedical Semantics 2018 9:12

### **Overall ECMT process**



**Table 7.** System performance for ICD10 coding on the **French raw** test corpus in terms of Precision (P), recall (R) and F-measure (F). A horizontal dash line places the frequency baseline performance. The top part of the table displays official runs, while the bottom part displays non-official and baseline runs.

ALL		EXTERNAL					
Team	P	R	$\mathbf{F}$	Team	P	R	$\mathbf{F}$
SIBM-run1	.857	.689	.764	SIBM-run1	.567	.431	.490
g LITL-run2	.666	.414	.510	LIRMM-run1	.443	.367	.401
$\mathbf{I}$ LIRMM-run1	.541	.480	.509	LIRMM-run2	.443	.367	.401
E LIRMM-run2	.540	.480	.508	LITL-run2	.560	.283	.376
E LITL-run1	.651	.404	.499	LITL-run1	.538	.277	.365
TUC-MI-run2	$.0\bar{4}4$	$.\bar{0}2\bar{6}$	$\overline{.033}$	$\overline{\mathrm{TUC}}$ - $\overline{\mathrm{MI}}$ -run $2$	$\overline{.010}$	004	$\overline{.005}$
TUC-MI-run1	.025	.015	.019	TUC-MI-run1	.006	.005	.005
average	.475	.358	.406	average	.367	.247	.292
median	.541	.414	.508	median	.443	.283	.376
LIMSI-run2	.872	.784	.825	LIMSI-run2	.700	.594	.643
E LIMSI-run1	.883	.760	.817	LIMSI-run1	.709	.559	.625
TUC-MI-run1-corrected	.883	.539	.669	TUC-MI-run1-corrected	.780	.290	.423
<b>5</b> TUC-MI-run2-corrected	.882	.536	.667	TUC-MI-run2-corrected	.767	.283	.414
uNIPD-run1	.629	.468	.537	UNIPD-run2	.350	.381	.365
${f \check{Z}}$ UNIPD-run2	.518	.384	.441	UNIPD-run1	.362	.251	.296
Mondeca-run1	$\overline{.375}$	.131	$.\bar{1}\bar{94}$	Mondeca-run1	.335	.228	.271
Frequency baseline	.339	.237	.279	Frequency baseline	.381	.110	.170

### Multiterminology Multilingual Semantic search engine

### Layer 3

- First step, definition of a model as light and as compact as possible
- Search engine based on Semantic Web
  - Explosion (subsumption) based on hierarchy at a multiterminology level
  - Other relations may be used: semantic expansions based on inter-terminology semantic mappings
  - Multilingual +++

Lelong, Romain; Soualmia, Lina F; Grosjean, Julien; Taalba, Mehdi & Darmoni, SJ. Building a Semantic Health Data Warehouse: Evaluation of a search tool in Clinical trials. JMIR Medical Informatics, 2019 (in press).

### Multiterminology Multilingual Semantic search engine

### Layer 3

- Very generic tool
- Designed to be used for N patients (regular use of HDW) or for one patient (not politically correct)
- use of HDW to consult patient data for CARE +++

Two or three generations of software (LIS, RIS, EHR) compiled in one HDW

Show me all the electroretinography reports for this patient

one second vs. at least one minute

# **TECHNICAL ASPECTS**



- Custom NoSQL (In Memory Data Grid IMDG-): Key-Value store, join indexes (replacing SQL joins) & Lucene (NLP)
- Three powerful servers:
  - 1 To RAM & 196 cores (prod)
  - 1 To RAM & 144 cores (preprod) \* 2
  - Semantic Annotator + deep learning









### Results

# **HDW Rouen - volumetry**

- 1.86 M patients
- 13.28 M stays (H C)
- 13.4 M documents (since 2000)
   5 G extracted medical concepts (ECMT)
- 121.1 M unitary biology tests (Na, K) (since 2004)
- DRG: 9.4 M diagnoses; 8.6 M procedures

### HDW Rouen – functional coverage





Soon available in 2018

# Main steps are already performed...

- A team of FIVE engineers are devoted to this project from DBMI + physicians from the DRG
- Scalability of the French semantic annotator
  - 11.8 M reports extracted in less then 24h; possible to rerun again every week if necessary (6ms for one report)
- Scalability of the semantic search engine
  - Using a NoSQL architecture and two powerful servers, response time ≈ 2 s
  - Might be used in care context
    - Display the last US of this patient ASAP

# Table 1 - List of different terminologies. with the number of annotations (before & after filtering)

Terminologies (SOC)	Number of annotations (after filtering)	Number of annotations (before filtering)	Filtering ratio
SNOMED CT	394 133 994	881 884 314	2.2
NCIt	319 853 952	843 195 067	2.6
MeSH	295 537 298	1 024 585 229	3.5
SNOMED 3.5	219 706 745	440 228 408	2.0
French Public Health Thesaurus	179 747 539	454 354 922	2.5
MedDRA	137 653 806	225 100 880	1.6
τυν	106 616 463	171 231 559	1.6
RADLEX	80 197 479	150 406 338	1.9
FMA	55 350 010	123 777 281	2.2
CISMeF	51 051 547	138 204 239	2.7

### Where ICD10 & CCA – French ICHI-?

### **Top 10 of Terminologies using unique annotated concepts (after filtering)**

Name of KOS	Total number of concepts	Number of translated concepts (%)	Number of unique identified concepts	Terminology coverage (%)
SNOMED CT Concept	326946	194611 (59.5)	<b>59330</b>	30.5
Notion SNOMED	100908	100908 (100)	36229	35.9
NCIt Concept	93925	68776 (73.2)	25315	36.8
MedDRA	44226	44226 (100)	22711	51.4
MeSH Descriptor	28329	28329 (100)	18288	64.6
MedDRA PT	21612	21612 (100)	13580	62.8
MeSH Concept	365731	102116 (27.9)	12625	12.4
FMA	81041	16629 (20.5)	8084	48.6
Radlex	42313	10259 (24.2)	6114	59.6
French Public Health	7097	7087 (100)	6080	85.0







#### Accès Sémantique à l'Information de Santé

Searched entity type :							
Patient Hospitalization	Stay	Patient managment	Biological test	Diagnosis	Procedure	Record	
Patient		- Gender	-	Mal	e Female	Other	
ET 👻 🕂 🕞 Diagnosis		- Terminolog	y(ies) -		• Enter	terms	
ET - + - Biological test		• Type of bio	logical test 🔹		Enter	terms	
ET 👻 🕂 🗖 Stay		• undefined	-		2018-04-	-30	
Construction de la requete :		Ajouter un	e ligne				
( @1 PATIENTS Male 105696 )							

Exemple de Requetes - Requête en language moteur

of concept provides Information Retrieval capabilities among a data subset of the health data warehouse of the university hospital of Rouen. Access to these data is regulated and an official request must be made at the

Health Data warehouse devoted to University Hospital of Rouen.

>













#### Quory buulding :

You can refer to above constraints. You can for instance type "@1" to refer to the first constraint or you can also use keywords such as "diagnosis", "patient", etc.

3.

2.

1,

Searched entity type :

Please select the types of entities that will be returned by the search engine

# Valorization

### Alicante

- SME from the North of France
- Exclusive sales of RUH DBI tools: HeTOP & ECMT => two French hospitals
- Integration of the semantic search engine in 2019

### OMICX

- One PhD to develop deep learning tools for annotation & information retrieval (sept. 2017 + sept. 2018)
- Licence fee for semantic search engine (2017-8)

# Use cases of the sHDW

- Since March 2018, over 60 use cases treated from SHDW (still POC)
  - Huge feedback to modify our top priorities to be developed
  - One example: to heavily invest on the semantic annotator as 80% of the questions were directly linked to health documents vs. lab tests or drugs
- One example: all lab tests during several years for specific infants
- In June 2018, to identify all the patients with endocarditis AFTER TAVI Transcatheter Aortic Valve Implantation
- To identify the TAVI acronym in the documents, thanks to the semantic annotator (ECMT)
  - Number of cases retrieved based on DRGs data warehouse = 30
  - Number of cases retrieved based on current RUH data warehouse = 53 (2 false negatives)
  - 23 new cases retrieved thanks to EDSaN!!!

# And soon...

- « v1 » to be launched (2019-2020)
  - Industrialization (Alicante)
  - « generic & unified » patient centric model +++
  - More data (drugs, medical devices, resuscitation, etc.)
  - Management of temporal data
    - Allen algebra (BEFORE, DURING, AFTER)
  - Query library
  - Export module for DRG department (SQL)
- HDW not limited to RUH: data from other hospitals, private offices, national data
- Health Data Hub in France ≠ Israeli (or Swiss) HDW

# **Doc'EDS**

- New tool developed in 2019, based on the first 20 use cases since March 2018
- No semantics... for now => better cost-effectiveness
- Based on Lucene
- To better understand the deeper structure of the document (80% of the health information in France)
  - Negation
  - Hypothesis
  - Family history
  - Segmentation of the document

# **Doc'EDS**

- Allow « manual dataminig »: oxymoron!
- Tool to validate the quality of the queries: continuous improvment of the search engine
- Based on rules, regular expressions and frequency analyses

	Query Part		Document part
Doc'EDS v0.1 - Re	ecense 13457785 documents de 2000 a	à février 201	9 Texte Méta-données Indevation Automatique Texte brut
Text doc.*	"tumeur de l'oesophage"	<b>\$</b> -	indexation Automation Automation Automation Medicales
Text doc. (négation)*		auto-completic	SCHEMA DE TRAITEMENT : LV5 FU2 - CISPLATINE 60 %.
Text doc. (condition)*		auto-completi(	CIBLES MESURABLES : Turneur primitive : oesophage.
Text doc. (atcd familiaux)*		auto-completic	
Date doc.			Asthénie ++
Type doc.*		auto-completic	EXAMEN CLINIQUE :
Unité(s) médicale(s)*		auto-completic	Indice de performance : grade OMS 1. Poids : 77 kg. Surf. corp.: 1.87 m2. Examen normal
UF(s)*		auto-completic	
DdN patient			EXAMENS COMPLEMENTAIRES :
Age (au moment du CR)			Biologie : GB : 2900 G/L PN : 2300 G/L Plaquettes : 104 G/L Hb : 8.8 mmol/L Hte : 25 % .
Sexe patient	1, 2		
Code(s) acte(s)*		auto-completic	CONCLUSION : Report de la cure LV5 FU2 - CISPLATINE 60 % pour carcinome épidermoïde du 1/3 supérieur et moyen de l'oesophage en raison d'une hématoxicité
Code(s) diag(s)*		auto-completic	grade II.
NIP			CONDUITE [DOCTOR] : Surveillance par NFS, urée et créatinine en externe.
ID CPAGE			Examen(s) programmé(s) : TDM le 01/12/05. Consultation [DOCTOR] le 07/12/05.
ID DOC			ł
Rechercher 1878 document(s) Exporter les résultats	1 de 1878 documer Précédent ▲ 001096850637 Partie « résu	ts > Suivant	B REHIMAT [DOCTOR]. DI FIORE Interne. // Courrier adressé à Madame le [DOCTOR] (DIEPPE), Monsieur le [DOCTOR] (LONGUEVILLE SUR SCIE), Monsieur le [DOCTOR] (ROUEN), Monsieur le [DOCTOR] (DIEPPE). // [LASTNAME] [FIRSTNAME]



FUROSEMIDE 20 mg, un par jour SELOKEN 200 mg, 1 par jour SINVASTATINE 40 mg, un par jour INNOHEP 14000, une injection par jour IMOVANE 7.5 mg, un au coucher VESICARE 5 mg, un par jour UROREC 8 mg, un par jour

MODE DE VIE Marié (femme aidante)

HISTOIRE DE LA MALADIE :

Patient de 
ans hospitalisé le matin 4/04 à Becquerel pour une cure de radiothérapie dans le caure de la prise en charge d'un carcinome épidermoide de l'oesophage suivit au CHU en gastro-entérologie (médecin référent : [DOCTOR]).

Le patient décrit des expectorations mêlées de sang survenues en deuxième moitié de nuit (nuit du ) au (04). Un second épisode de crachats sanglants après effort de toux au cours de l'évaluation pré radiothérapie. Dans le contexte de tumeur de l'oesophage, le patient est transféré en unité de soins intensifs de gastro-entérologie au CHU pour une probable hématémèse.

EXAMEN CLINIQUE :

Taille : 176 cm : 76 kg. IMC : 24.5 kg/m². Température : 36.4°C Pouls : 77 par mn. TA : 109/70. EVA : 0 de l'hémorragie, pas de déglobulisation

Dénutrition. OMS 2-3.

Conjonctives colorées. Absence d'ictère.

Cardiovasculaire : BDC irrégulières. Pas de souffle. Pas d'OMI Pas de Turgescence des jugulaires.

Respiratoire : Râles crépitants à la base pulmonaire gauche. Toux importantes suivies d'expectorations teintées de sang. Digestif : présence d'une sonde naso-gastrique d'alimentation entérale. Abdomen souple, indolore. BHA présents et normaux. Absence d'argument en faveur d'une ascite clinique. Toucher rectal : selles dures de couleur normale (pas de méléna ni de rectorragies). Prostate de taille augmentée avec disparition du sillon médian.

Neurologique : bonne orientation temporo-spatiale. Pas de déficit sensitivo-moteur.

BIOLOGIE :

Hémoglobine à 12,3 g/dl. .Sérologie et antigénurie aspergillaire en cours.

RADIOLOGIE :

Radio de thorax :

TDM du 2/20 : Stabilité de l'épaississement circonférentiel oesophagien avec multiples adénomégalies médiastinales. Apparition de plusieurs nodules lobaires inférieures gauches :

### Detection of

Anonymization

(patients and

physicians)

### suspiscion/hypotheses/

doubts/future

### Détection of negations



### Différents bilans

### Analyses

#### Statistiques & bilans



150 femmes (22,4%)	
520 hommes (77,6%)	
Statistiques des ages (au moment du CR)	
Moyenne	63,3
Écart type	11,5
Minimum	7
Maximum	93
Q1	54
Q2	63
Q3	72







		-	-
DIGE HEPATO GASTRO ENTEROLOGIE NUTRITION		1685	61,9%
CGCD CHIRURGIE GENERALE ET DIGESTIVE		235	8,6%
RADI IMAGERIE CENTRALE		149	5,5%
URGE URGENCES	_	128	4,7%
PHIE PHARMACIE		93	3,4%
ORLO O.R.L ADULTES	Répartition dans les	61	2,2%
PNM1 CLINIQUE PNEUMOLOGIQUE HCN	unités médicales /	59	2,2%
PHYS PHYSIOLOGIE DIGESTIVE		34	1,2%
CARD CARDIOLOGIE			0,9%
REAC REANIMATION CHIRURGICALE		23	0,8%

### Diagnostics PMSI, etc.

Afficher 10 - éléments

Code(s) diag(s)

DA DR

Filtrer :

155 valeurs distinctes, 2369 occurrences au total

DP

DIAGCODES		\$
Z511 séance de chimiothérapie pour tumeur	1292 54,5%	
Z530 acte non effectué en raison de contre-indication	302 12,7%	
C155 tumeur maligne du tiers inférieur de l'oesophage	86 3,6%	
C151 tumeur maligne de l'oesophage thoracique	66 2,8%	
Z087 examen de contrôle après traitements combinés pour tumeur maligne	63 2,7%	
C154 tumeur maligne du tiers moyen de l'oesophage	50 2,1%	
C150 tumeur maligne de l'oesophage cervical	39 1,6%	
C153 tumeur maligne du tiers supérieur de l'oesophage	39 1,6%	
Z452 ajustement et entretien d'un dispositif d'accès vasculaire	34 1,4%	
C159 tumeur maligne de l'oesophage, sans précision	24 1,0%	
Affichage de l'élement 1 à 10 sur 155 éléments		
	Précédent         1         2         3         4         5          16	5

Suivant

# Data mining (ECMT)

Indexations automatiques couvrant plus	s d'un document			×
Afficher	l concepts ected corp	in us	Filtrer :	Ł
Catégorie	Concept	dentifiant	Terminologie	♦ Nb documents ♦
processus néoplasique;	tumeur	SCT_CO_108369006	SCT	2313
thérapeutique; médicaments; pharmacie; procedure thérapeutique ou préventive;	chimiothérapie	MSH_D_004358	MSH	1856
procedure thérapeutique ou préventive;	chimiothérapie	SCT_CO_363688001	SCT	1856
procedure thérapeutique ou préventive;	chimiothérapie	SCT_CO_367336001	SCT	1856
diagnostic; activité de soins médicaux;	examen clinique	MSH_D_010808	MSH	1801
<b>EXAM</b> <b>EXAM</b> <b>INTROLIT TADIOTHÉRA</b>	en cli Carc dm	Précédent 1	2 3 742	4 5 Suivant



Affichage de l'élement 1 à 5 sur 332 éléments (filtré de 3,710 éléments au total)



# Perspectives

- Medical word embeddings (Word2Vec, FastText,Glove)
  - PhD Emeric Dynomant (OMICX)
  - Two different corpus
    - 12 M health documents from HSDW
    - 180 K abstracts from LiSSa, French bibliographic database
- Doc2Vec, Patient2Vec (in progress)
  - PhD Mikaël Dusenne, MD
    - Hybrid semantic annotator ("old" NLP + deep NLP)
    - Doc2Vec2DRGs, using an other tool !!! ELMO?



12M • endocardite	Search
GloVe	infectieuse, myocardite, eto, native, streptocoque, bovis, bactériémie, faecalis, _endocardite
FastText (CBOW)	proprio_septive, myopericardite, septo_optique, endo_aortique, endocartite, endoculaire, acrodermite, rhino_septale, salmonellose, épidermolyse
FastText (Skip-Gram)	endocartite, endocardique, proctologique, extancilline, septo_basale, prolongements, recanalisée, précentrale, dantrolene, podoscopique
Word2Vec (Skip-Gram)	bovis, sanguinis, eto, gordonii, gallolyticus, aorto_mitrale, mutans, infectieuse, streptoccoque, salivarius
Word2Vec (CBOW)	endocartite, _endocardite, native, bovis, médiastinite, myocardite, mutans, gallolyticus, myopéricardite, tamponnade

Home

ervation Cinémas Pathé × 🗅 Medical Embedding - Result × 😒 darma	ni - PubMed - NCBI 🗙 📘 Clinical Natural Language Pr 🗙 🥰 LIMICS	🗙 🛛 🗾 licence pour logiciel - Tradu 🗙 🕇 💾 H
	Medical word embeddin	ngs querying page
endocardite Search		
ord2Vec (Skip-Gram) endocardites, bactériémie, infectieuse,	valvulopathie, septicémie, mycotique, spondylodiscite, médiastinite, v	alvulaire, bioprothèse

Home	

# Wordembeddings in two different contexts

QUERY: "facebook"

LiSSa corpus (300k)	internet, twitter, web, blog, e_learning, blogs, internautes, tic, game,
RUH documents (12M)	reproches, injures, messages, insultes, rumeurs, ex_conjointe, menaces, insultant,

### Doc2Vec2PubMed



Thank you for listening!

**Questions?** 

Email: Stefan.Darmoni@chu-rouen.fr