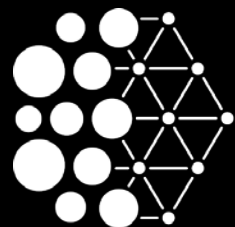


Institut  
québécois  
d'intelligence  
artificielle



Mila

## Applications : Extraction de texte

Joseph Paul Cohen, PhD (joseph@josephpcohen.com)  
Margaux Luck, PhD (margaux.luck@mila.quebec)

1er décembre 2018

# Extraction de texte

Revue par Joseph Paul Cohen, PhD  
Mila - Institut Québécois d'Intelligence Artificielle

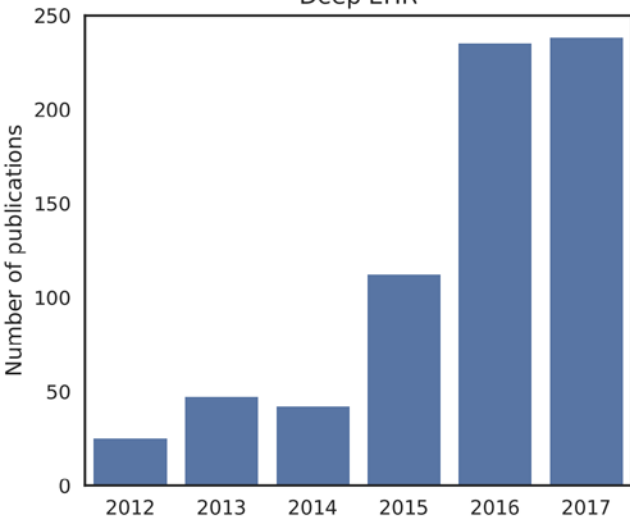
Sujets:

1. Représentation de concepts médicaux
2. Word2Vec

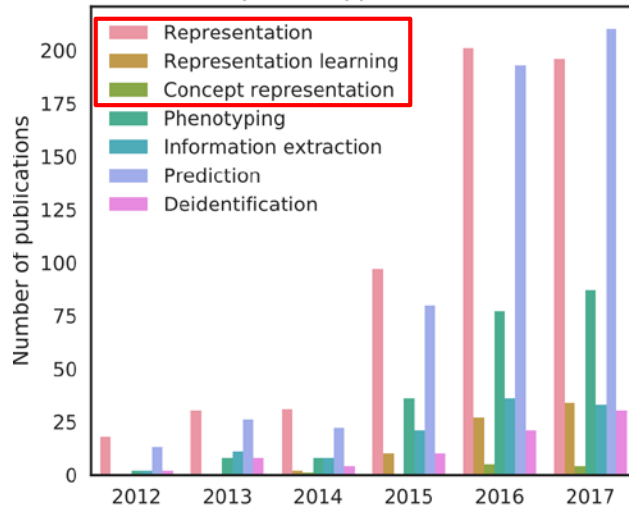


# Où en est la recherche en apprentissage profond ?

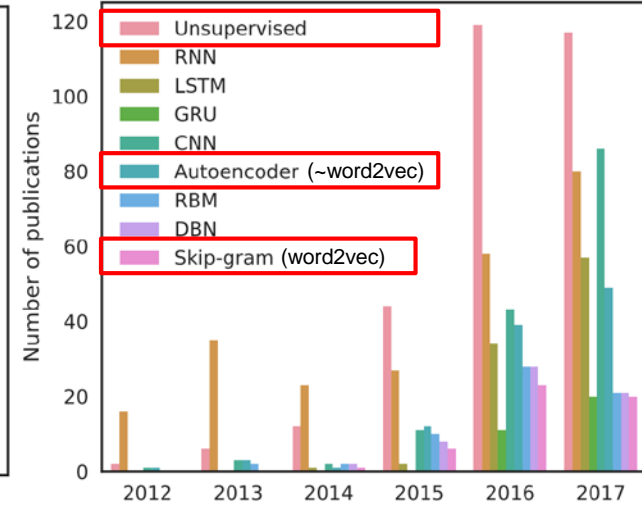
Deep EHR



Deep EHR: Application Areas



Deep EHR: Technical Methods



## Deep EHR: A Survey of Recent Advances in Deep Learning Techniques for Electronic Health Record (EHR) Analysis

Benjamin Shickel<sup>1</sup>, Patrick J. Tighe<sup>2</sup>, Alex Himeva<sup>3</sup>, and Robin Rindt<sup>1</sup>

Abstract—The past decade has seen an explosion in the amount of digital information stored in electronic health records (EHR). This paper is designed for academic and professional administrators, health-care providers, and researchers interested in the use of deep learning for EHR analysis. We survey the state-of-the-art in deep learning for EHR analysis, covering a wide range of applications including supervised classification, representation learning, and unsupervised representation learning. We discuss the benefits and challenges of deep learning for EHR analysis, and provide a list of key research papers and resources for further reading.

[Shickel, Deep EHR : A Survey of Recent Advances in Deep Learning Techniques for Electronic Health Record, 2018]

# Représentation de concepts

## Note clinique

Mr. Smith is a 63-year-old gentleman with coronary artery disease, hypertension, hypercholesterolemia, COPD and tobacco abuse. He reports doing well. He did have some more knee pain for a few weeks, but this has resolved. He is having more trouble with his sinuses. I had started him on Flonase back in December. He says this has not really helped. Over the past couple weeks he has had significant congestion and thick discharge. No fevers or headaches but does have diffuse upper right-sided teeth. **He denies any chest pains, palpitations, PND, orthopnea, and syncope.** His breathing is doing No cough. He continues to smoke half-a-pack per day. He plans on **trying the patches again.**

## Publication clinique

The screenshot shows the NCBI PubMed interface for a clinical publication. The page title is "Collaborative Efforts Driving Progress in Pediatric Acute Myeloid Leukemia". The authors listed are C. Michel Zwaan, Edward A. Kolb, Dirk Reinhardt, Jonas Abrahamsson, Souichi Adachi, Richard Aclero, Eveline S.-J.M. De Bont, Barbara De Moerloose, Michael Dworzak, Brenda E.S. Gibson, Henrik Hasle, Guy Levisser, Franco Locatelli, Christine Ragu, Raul C. Ribeiro, Carmelo Rizzari, Jeffrey E. Rubnitz, Owen P. Smith, Lillian Sung, Daisuke Tomizawa, Marry M. van den Heuvel-Eibrink, Ursula Creutzig, and Gerrit J.L. Kaspers. The publication date is 2015 Aug 24. The PMID is 26304895. The page includes a search bar and navigation links like "Home", "This Article", "Search", and "Submit".

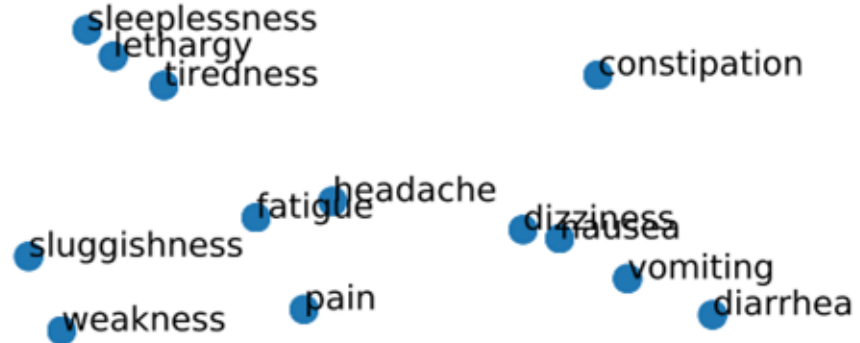
Utile pour comprendre les similitudes ou faire des prédictions semi-supervisées !

## Représentations pour:

- Les patients
- Les médecins
- Les visites
- Les maladies
- Les médicaments
- Les symptômes

# Encodage de mots dans le langage biomédical

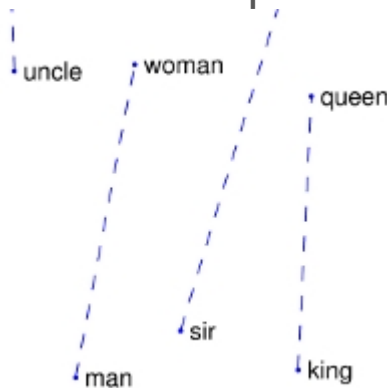
Extraire les relations  
entre les mots et produire  
une représentation  
latente



# Que faire avec l'encodage de mots?

- Nous pouvons les composer pour créer des encodages de paragraphe.
- Utilisation à la place des mots pour un RNNs
- Augmenter les représentations apprises sur de petits ensembles de données

Étudier la compositionnalité de l'espace latent appris

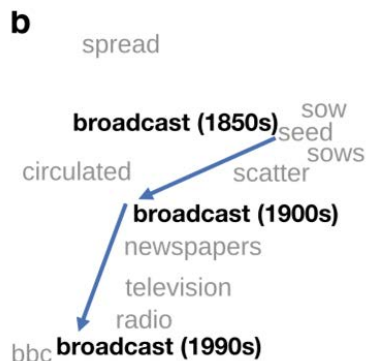


[Mikolov, 2013]

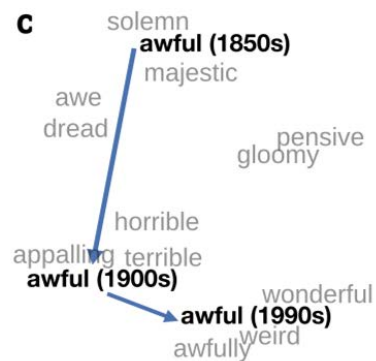


[Pennington, 2014]

Étudiez la variation de sens entre deux textes (ou hôpitaux, ou médecins) ?



[Cultural Shift or Linguistic Drift, Hamilton, 2016]



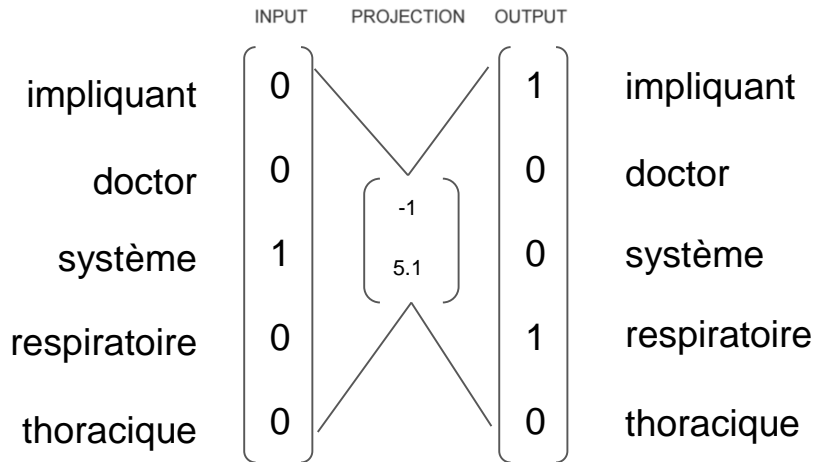


# word2vec

contexte fenêtre



1. Chaque mot est un exemple d'entraînement
2. Chaque mot est utilisé dans de nombreux contextes
3. Le contexte définit chaque mot



## Efficient Estimation of Word Representations in Vector Space

Tomas Mikolov  
Google Inc., Mountain View, CA  
tmikolov@google.com

Kai Chen  
Google Inc., Mountain View, CA  
kaichen@google.com

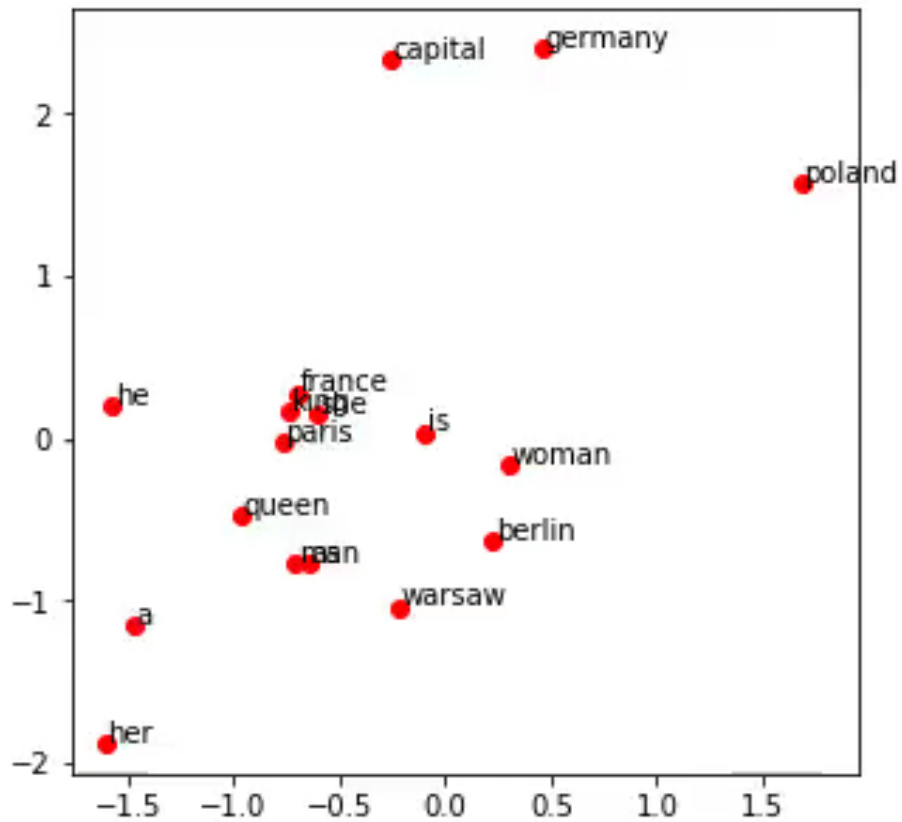
Greg Corrado  
Google Inc., Mountain View, CA  
gcorrado@google.com

Jeffrey Dean  
Google Inc., Mountain View, CA  
jdf@google.com

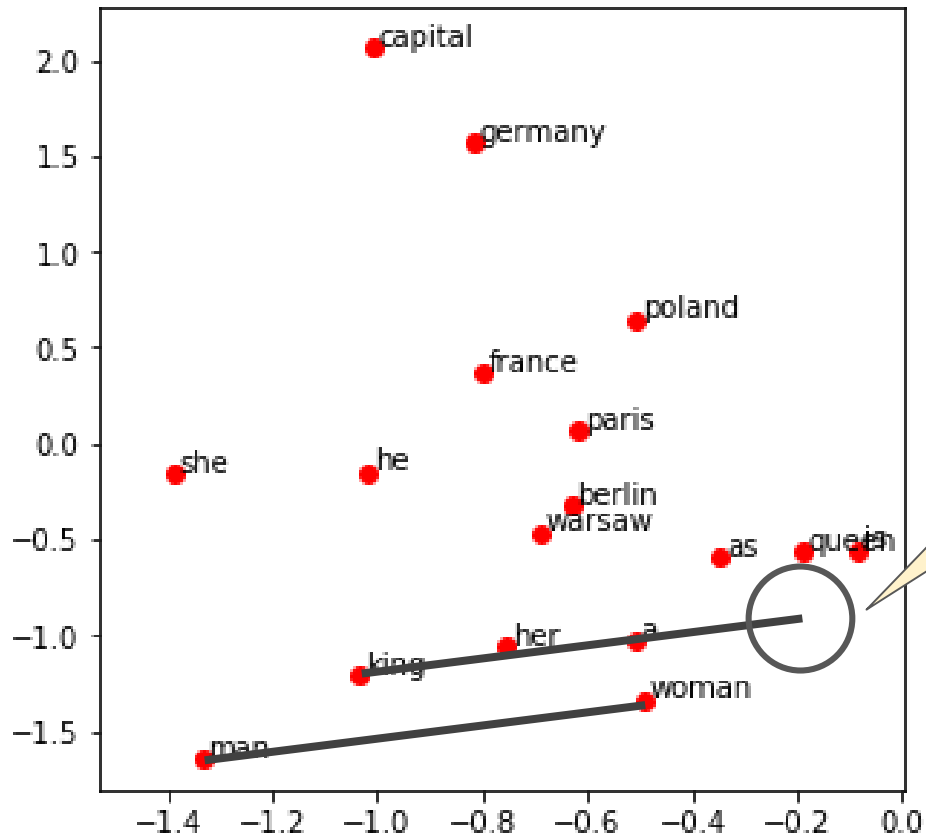
### Abstract

We propose two novel neural architectures for computing continuous vector representations of words from very large data sets. The quality of these representations is measured in a word similarity task, and the results are compared to the state-of-the-art neural approaches. We observe that the quality of the word representations is sensitive to the choice of neural network architecture, and that the quality of the word representations is sensitive to the choice of neural network architecture. We observe that the quality of the word representations is sensitive to the choice of neural network architecture, and that the quality of the word representations is sensitive to the choice of neural network architecture.





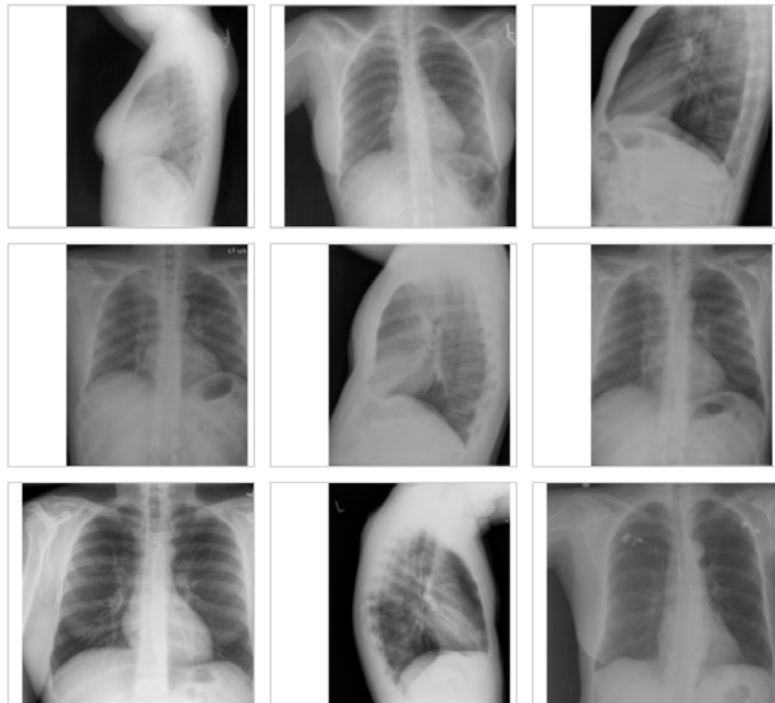
Apprentissage en cours



roi + (femme - homme) = ?

# Rapports de l'hôpital universitaire de l'Indiana

Images de radiographie pulmonaire du réseau hospitalier de l'Université de l'Indiana



## Indiana University Chest X-ray Collection

Kohli MD, Rosenman M - (2013)

**Affiliation:** Indiana University

### ABSTRACT

**Comparison:** None.

**Indication:** Positive TB test

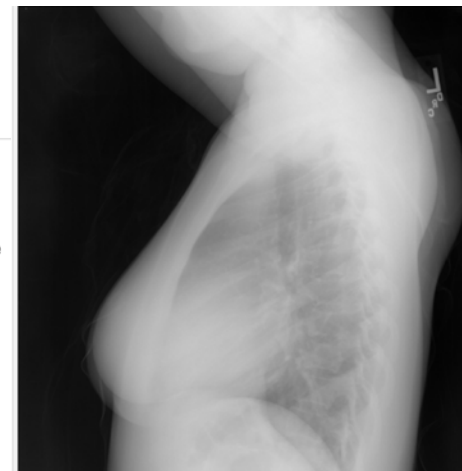
**Findings:** The cardiac silhouette and mediastinum size are within normal limits. There is no pulmonary edema. There is no focal consolidation. There are no XXXX of a pleural effusion. There is no evidence of pneumothorax.

**Impression:** Normal chest x-XXXX.

**NOTE:** The data are drawn from multiple hospital systems.

[Show MeSH](#)

Related in: [MedlinePlus Request Collection](#)



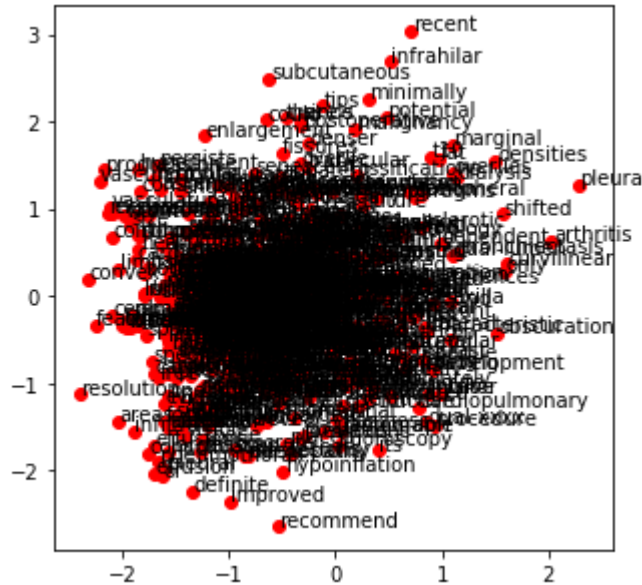
4000 rapports disponibles en format XML !

OPEN 

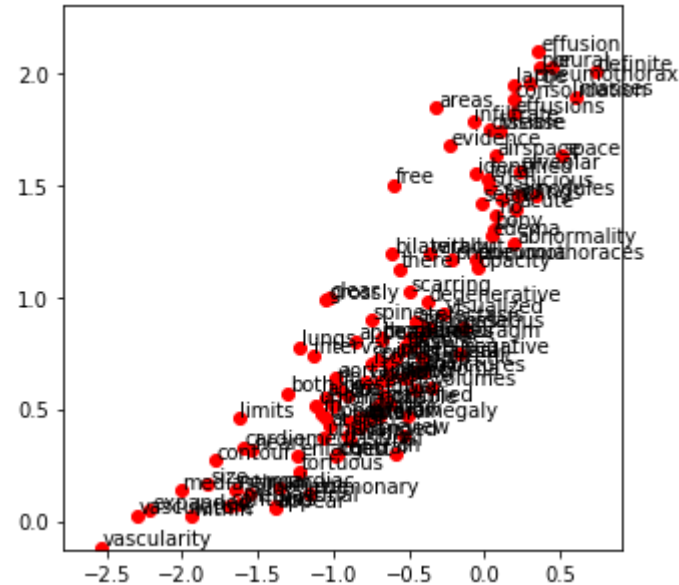
# Hyperparamètres !

En fonction de la configuration du modèle, les encodages peuvent varier :

Aucune information



Compression



Nous pouvons varier :

la dimension de l'encodage, le taux d'apprentissage, les mots-clés, la taille de la fenêtre, etc.



## Sous-ensemble de fichiers disponibles en libre accès

- Documents disponibles en format PDF ou XML
- 1,25 million d'articles biomédicaux et 2 millions de mots distincts
- Disponible par FTP pour téléchargement en masse
- Les métadonnées comprennent le nom de la revue et l'année.

```
</sec>
<sec sec-type="methods">
  <title>Patients and methods</title>
  <sec>
    <title>Patients</title>
    <p>Between September 1995 and September 1998, 413 patients with abnormal breast findings were referred for histological evaluation to the Department of Gynecology of the Friedrich-Schiller University, Jena, Germany. Patients had been selected and referred because of the presence of breast lesions detected by palpation and/or mammography and/or sonography. In addition, MR mammography was performed in all patients. We excluded five patients with invasive cancer who had a history of core-needle or fine-needle biopsy cancer within 2 weeks before referral, because the presence of haematoma may mimic false-positive findings on MR mammography. In addition, five patients who did not keep still during MR mammography were excluded.</p>
  </sec>
  <sec>
    <title>Imaging</title>
    <p>
      Analysis of the sonograms taken in patients with histologically confirmed carcinoma
      <itali>in situ</itali>
      were excluded from analysis because the value of sonography for detection of premalignant disease is
```

Exemple de données XML

```
Breast Cancer Res.
Genome Biol.
Arthritis Res.
BMC Cell Biol.
...
```

Noms des journaux

<https://www.ncbi.nlm.nih.gov/pmc/tools/openftlist/>



Exemples de mots médicaux et leurs cinq mots similaires basés sur chaque corpus de texte de formation

Le nom complet du diabète est "diabète mellitus"

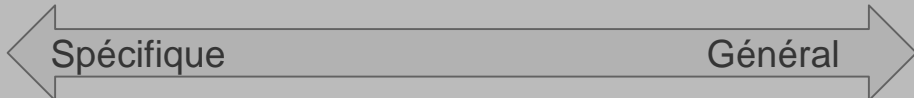
Selon certains articles, les diabétiques courent un risque accru d'hypertension (hypertension artérielle)

Un type d'ulcère gastro-duodéal

Un des symptômes est un ulcère

Le cancer du côlon est associé au cancer du sein !

	EHR (Mayo Clinic)	PubMed	Wikipedia + Gigaword	Google N
diabetes	mellitus,	cardiovascular,	hypertension,	diabetics,
	controlled,	nonalcoholic,	obesity,	hypertension,
	hypersterolemia,	obesity,	arthritis,	diabetic,
	epidemia,	mellitus,	cancer,	diabetes_mellitus,
	diabetes_mellitus	polycystic	alzheimer	heart_disease
peptic ulcer disease	scleroderma,	gastritis,	ulcers,	ichen_planus,
	duodenal,	alcoholism,	arthritis,	Candida_infection,
	crohn,	rheumatic,	diseases,	vaginal_yeast_infections,
	gastroduodenal,	ischaemic,	diabetes,	oral_thrush,
	diverticular	nephropathy	stomach	dermopathy
colon cancer	breast,	breast,	breast,	breast,
	ovarian,	mcf,	prostate,	prostate,
	prostate,	cancers,	cancers,	tumor,
	postmenopausally,	tumor_suppressing,	tumor,	pre_cancerous_lesion,
	caner	downregulation	liver	cancerous_polyp







Joseph Paul Cohen,  
PhD



Myriam Côté,  
PhD



Pr. Yoshua Bengio,  
PhD



Mandana Samiei



Francis Dutil



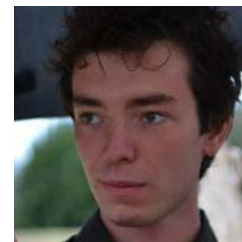
Martin Weiss



Shawn Tan



Geneviève Boucher



Georgy Derevyanko,  
PhD



Tristan Sylvain



Margaux Luck,  
PhD



Sina Honari



Assya Trofimov



Vincent Frappier,  
PhD

# Merci

## Références

Ching, T., et al. **Opportunities And Obstacles For Deep Learning In Biology And Medicine**. Journal of The Royal Society Interface. 2018

Shickel, B. et al. **Deep EHR : A Survey of Recent Advances in Deep Learning Techniques for Electronic Health Record**. IEEE Journal of Biomedical and Health Informatics, 2018