



# AfIA

Association française  
pour l'Intelligence Artificielle

## DEFT

---

*Défi Fouille de Textes  
(atelier TALN-RECITAL)*

---

## PFIA 2019





## Table des matières

Natalia Grabar, Cyril Grouin, Thierry Hamon et Vincent Claveau. <b>Éditorial</b> .....	4
. <b>Comités</b> .....	5
Natalia Grabar, Cyril Grouin, Thierry Hamon et Vincent Claveau. <b>Recherche et extraction d'information dans des cas cliniques. Présentation de la campagne d'évaluation DEFT 2019</b> .....	7
Philippe Suignard, Meryl Bothua et Alexandra Benamar. <b>Participation d'EDF R&amp;D à DEFT 2019 : des vecteurs et des règles!</b> .....	17
Jacques Hilbey, Louise Deléger et Xavier Tannier. <b>Participation de l'équipe LAI à DEFT 2019</b> .....	29
Jean-Christophe Mensonides, Pierre-Antoine Jean, Andon Tchechmedjiev et Sébastien Harispe. <b>DÉfi Fouille de Textes 2019 : indexation par extraction et appariement textuel</b> .....	35
Davide Buscaldi, Dhaou Ghoul, Joseph Le Roux et Gaël Lejeune. <b>Indexation et appariements de documents cliniques pour le Deft 2019</b> .....	49
Mérimèe Bouhandi, Florian Boudin et Ygor Gallina. <b>DeFT 2019 : Auto-encodeurs, Gradient Boosting et combinaisons de modèles pour l'identification automatique de mots-clés. Participation de l'équipe TALN du LS2N</b> .....	57
Estelle Maudet, Oralie Cattan, Maureen de Seyssel et Christophe Servan. <b>Qwant Research @DEFT 2019 : appariement de documents et extraction d'informations à partir de cas cliniques</b> .....	67
Damien Sileo, Tim Van de Cruys, Philippe Muller et Camille Pradel. <b>Apprentissage non-supervisé pour l'appariement et l'étiquetage de cas cliniques en français - DEFT2019</b> .....	81
Khadim Dramé, Ibrahima Diop, Lamine Faty et Birame Ndoeye. <b>Indexation et appariement de documents cliniques avec le modèle vectoriel</b> .....	91

# Éditorial

La reproductibilité des résultats et la robustesse des outils constituent des enjeux critiques en traitement automatique des langues, en particulier parce que le TAL commence à fournir des méthodes et outils mûrs qui sont de plus en plus utilisés dans d'autres domaines, comme le domaine médical. L'objectif des compétitions en TAL est de fournir des corpus et les données de référence qui permettent aux chercheurs de développer des outils et de les tester ensuite. Un tel contexte permet également d'avoir une première comparaison entre les méthodes et approches utilisées par les participants du défi, dans des conditions expérimentales parfaitement identiques.

L'édition 2019 du défi fouille de textes (DEFT 2019, <https://deft.limsi.fr/2019/>) a porté sur l'analyse de cas cliniques rédigés en français. Trois tâches ont été proposées autour de la recherche d'information et de l'extraction d'information, en s'inspirant de tâches réelles et utiles pour le domaine médical. La particularité de cette édition concerne ainsi le domaine traité (médical) et les documents utilisés (des cas cliniques). C'est la première fois qu'une compétition a lieu sur des textes cliniques en français. Les cas cliniques décrivent les situations cliniques de patients, réels ou fictifs. Ils sont publiés dans différentes sources de données (scientifique, didactique, associatif, juridique, etc.), de manière anonymisée. L'utilité des cas consiste à présenter des situations cliniques typiques ou rares, notamment à des fins pédagogiques. Le corpus de cas cliniques utilisé lors de la campagne DEFT 2019 se compose de cas librement accessibles en ligne.

Lors du déroulement de la campagne, l'accès aux données d'entraînement a été possible dès le 18 février, tandis que la phase de test s'est déroulée du 9 au 15 mai, sur une période de trois jours définie par chacun des participants. Huit équipes se sont inscrites et ont participé jusqu'au bout. Nous comptons cinq équipes académiques (LGI2P/Mines Alès, Nîmes; LIMICS/INRA, Paris; LIPN/STIH, Paris; TALN-LS2N, Nantes; Université Assane Seck de Ziguinchor, Sénégal), deux équipes industrielles (EDF Lab, Palaiseau; Qwant, Paris) et une équipe mixte (Synapse/IRIT, Toulouse).

Ces actes rassemblent la présentation des objectifs de la campagne (corpus, tâches, évaluation...), les résultats obtenus sur les différentes tâches et la description des systèmes participants.

Les organisateurs remercient le comité de programme pour avoir apporté leur soutien et leur expertise à la campagne d'évaluation DEFT 2019.

Natalia Grabar, Cyril Grouin, Thierry Hamon et Vincent Claveau

# Comités

## Organisateurs de DEFT

- Vincent CLAVEAU (IRISA, CNRS)
- Natalia GRABAR (STL, CNRS, Université de Lille)
- Cyril GROUIN (LIMSI, CNRS, Université Paris-Saclay)
- Thierry HAMON (LIMSI, CNRS, Université Paris-Saclay ; Université Paris XIII)

## Comité de programme de DEFT

- Patrice BELLOT (LSIS, Aix-Marseille Université)
- Leonardo CAMPILLOS LLANOS (LIMSI, CNRS, Université Paris-Saclay ; Madrid)
- Vincent CLAVEAU (IRISA, CNRS)
- Natalia GRABAR (STL, CNRS, Université de Lille)
- Cyril GROUIN (LIMSI, CNRS, Université Paris-Saclay)
- Vincent GUIGUE (LIP6, Sorbonne Université)
- Thierry HAMON (LIMSI, CNRS, Université Paris-Saclay ; Université Paris XIII)
- Véronique MORICEAU (LIMSI, Université Paris-Sud, Université Paris-Saclay ; IRIT)
- Fleur MOUGIN (Bordeaux Population Health, Université de Bordeaux)
- Mathieu ROCHE (TETIS, CIRAD)
- Patrick RUCH (HEG Geneva, BiTeM)
- Frantz THIESSARD (Bordeaux Population Health, Université de Bordeaux, Inserm ; CHU de Bordeaux, SIM pôle santé publique, unité médicale Informatique et archivistique médicales)

