# Pedro **Ortiz Suarez**
DOCTORANT

✉ pedro.ortiz@inria.fr | ⌂ https ://pjortiz.fr/ | ⬚ pjox

## **Int**érêts

Language Modeling, Sequence Tagging, Reconnaissance d'Entités Nommés, Parsing, Linguistique de Corpus, Text Mining.

## **For**mation

**Sorbonne Université**                                                                                            *Paris, France*
DOCTORAT EN INFORMATIQUE                                                                   *Oct. 2018 - Présent*
Contrat doctoral Inria

**Université de Paris 8, Vincennes–Saint-Denis**                                 *Saint-Denis, France*
LICENCE MATHÉMATIQUES, INFORMATIQUE APPLIQUÉES AUX SCIENCES HUMAINES ET SOCIALES     *Nov. 2017 - Sep. 2018*
Option sciences du langage

**Université d'Aix Marseille**                                                                         *Marseille, France*
MASTER EN MATHÉMATIQUES ET APPLICATIONS, PARCOURS MATHÉMATIQUES FONDAMENTALES     *Sep. 2016 - Jun. 2017*
Bourse d'excellence académique LabEx Archimède

**Universidad Nacional de Colombia**                                                         *Medellín, Colombie*
PREMIER CYCLE EN MATHÉMATIQUES FONDAMENTALES (BAC+4)                    *Jan. 2012 - Jun. 2016*
Admission honorable, troisième place au test d'admission

## **Exp**érience

**Inria - ALMAnaCH Team**                                                                           *Paris, France*
DOCTORANT                                                                                           *Nov. 2018 - Présent*
• Recherche en intelligence artificielle appliqué au traitement des langues naturelles.

**Sorbonne Université**                                                                               *Paris, France*
CHARGÉ DE TRAVAUX DIRIGÉS                                                               *Sep. 2019 - Jun. 2021*
• Mathématiques discrètes, TD 81h.
• Eléments de programmation 2, TME 30h.

**Inria - ALMAnaCH Team**                                                                           *Paris, France*
RELAIS THÈSE                                                                                               *Oct. 2018*
• Recherche en intelligence artificielle appliqué au traitement des langues naturelles.

**Inria - ALMAnaCH Team**                                                                           *Paris, France*
STAGIAIRE                                                                                           *Apr. 2018 - Sep. 2018*
• Recherche en text mining et intelligence artificielle appliqué aux accords d'entreprise.

**Institut de Mathématiques de Marseille, I2M**                                 *Marseille, France*
STAGIAIRE                                                                                           *Avr. 2017 - Jun. 2017*
• Recherche en géométrie algébrique complexe, en surfaces canoniquement plongées de grands degrés.

**Universidad Nacional de Colombia**                                                         *Medellín, Colombie*
CHARGÉ DE TRAVAUX DIRIGÉS                                                               *Aoû. 2015 - Jun. 2016*
• Mathématiques discrètes, TD 128h.

# Publications

[1] Julien Abadji, Pedro Javier Ortiz Suárez, Laurent Romary and Benoît Sagot. *Ungoliant : An Optimized Pipeline for the Generation of a Very Large-Scale Multilingual Web Corpus*, Proceedings of the Workshop on Challenges in the Management of Large Corpora (CMLC-9), Jul 2021, Leibniz-Institut für Deutsche Sprache, Online.

[2] Isaac Caswell, Julia Kreutzer, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmungkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Javier Ortiz Suárez, Iroro Orife, Kelechi Ogueji, Rubungo Andre Niyongabo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhalov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Ballı, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, Mofetoluwa Adeyemi. *Quality at a Glance : An Audit of Web-Crawled Multilingual Datasets*, arXiv e-prints, Mar 2021, arXiv, Online.

[3] Pedro Javier Ortiz Suárez, Yoann Dupont, Gaël Lejeune, Tian Tian. *SinNer@Clef-Hipe2020 : Sinful adaptation of SotA models for Named Entity Recognition in French and German*, Working Notes of CLEF 2020 - Conference and Labs of the Evaluation Forum, Sep 2020, CEUR-WS, Thessaloniki / Virtual, Greece.

[4] Pedro Javier Ortiz Suárez, Benoît Sagot, Laurent Romary. *A Monolingual Approach to Contextualized Word Embeddings for Mid-Resource Languages*, Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Jul 2020, Association for Computational Linguistics, Online.

[5] Louis Martin[*], Benjamin Muller[*], Pedro Javier Ortiz Suárez[*], Yoann Dupont, Laurent Romary, Éric de la Clergerie, Djamé Seddah, Benoît Sagot. *CamemBERT : a Tasty French Language Model*, Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Jul 2020, Association for Computational Linguistics, Online.

[6] Djamé Seddah, Farah Essaidi, Amal Fethi, Matthieu Futeral, Benjamin Muller, Pedro Javier Ortiz Suárez, Benoît Sagot, Abhishek Srivastava. *Building a User-Generated Content North-African Arabizi Treebank : Tackling Hell*, Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Jul 2020, Association for Computational Linguistics, Online.

[7] Louis Martin[*], Benjamin Muller[*], Pedro Javier Ortiz Suárez[*], Yoann Dupont, Laurent Romary, Éric de la Clergerie, Djamé Seddah, Benoît Sagot. *Les modèles de langue contextuels Camembert pour le Français : impact de la taille et de l'hétérogénéité des données d'entrainement*, 27e conférence sur le Traitement Automatique des Langues Naturelles (TALN 2020), Jun 2020, Nancy, France.

[8] Pedro Javier Ortiz Suárez, Yoann Dupont, Benjamin Muller, Laurent Romary, Benoît Sagot. *Establishing a New State-of-the-Art for French Named Entity Recognition*, 12th International Conference on Language Resources and Evaluation (LREC 2020), May 2020, Marseille, France.

[9] Murielle Popa-Fabre, Pedro Javier Ortiz Suárez, Benoît Sagot, Éric de la Clergerie. *French Contextualized Word-Embeddings with a sip of CaBeRnet : a New French Balanced Reference Corpus*, 8th Workshop on the Challenges in the Management of Large Corpora (CMLC-8), May 2020, Marseille, France.

[10] Mohamed Khemakhem, Ioana Galleron, Geoffrey Williams, Laurent Romary, Pedro Javier Ortiz Suárez. *How OCR Performance can Impact on the Automatic Extraction of Dictionary Content Structures*, 19th annual Conference and Members' Meeting of the Text Encoding Initiative Consortium (TEI) -What is text, really ? TEI and beyond, Sep 2019, Graz, Austria.

[11] Pedro Javier Ortiz Suárez, Benoît Sagot, Laurent Romary. *Asynchronous Pipelines for Processing Huge Corpora on Medium to Low Resource Infrastructures*, 7th Workshop on the Challenges in the Management of Large Corpora (CMLC-7), Jul 2019, Cardiff, United Kingdom.

[*]Contribué à parts égales.

## Talks

**Séminaires du Lattice** *Paris, France*
DES MÉTHODES DE TAL MODERNES POUR L'ENRICHISSEMENT DE DOCUMENTS *22/09/2020*
École normale supérieure / Laboratoire Lattice

**58ᵗʰ Annual Meeting of the Association for Computational Linguistics** *En Ligne*
A MONOLINGUAL APPROACH TO CONTEXTUALIZED WORD EMBEDDINGS FOR MID-RESOURCE LANGUAGES *06/07/2020*
Association for Computational Linguistics

**7ᵗʰ Workshop on the Challenges in the Management of Large Corpora** *Cardiff, Royaume-Uni*
ASYNCHRONOUS PIPELINE FOR PROCESSING HUGE CORPORA ON MEDIUM TO LOW RESOURCE INFRASTRUCTURES *22/07/2019*
Cardiff University

**10ᵗʰ International Conference on Historical Lexicography and Lexicology** *Leeuwarden, Pays-Bas*
PREPARING THE DICTIONNAIRE UNIVERSEL FOR AUTOMATIC ENRICHMENT *13/06/2019*
Fryske Akademy

**GIG #3 : bring the cool back in the cloud** *Paris, France*
REDUCING COMPUTATION TIME BY MONTHS BY REWRITING BASH SCRIPTS IN GO *24/03/2019*
Golang Paris

## Distinctions

2016-2017 **Attribution bourse d'excellence Master 2**, Laboratoire d'excellence LabEx Archimède. *Marseille, France*
Jan. 2012 **Bourse au mérite**, Universidad Nacional de Colombia. *Medellín, Colombie*
Aoû. 2012 **Bourse au mérite**, Universidad Nacional de Colombia. *Medellín, Colombie*

## Langues

**Espagnol** - Langue maternelle, **Anglais** - Courant, **Français** - Courant, **Allemand** - Élémentaire.

## Compétences informatiques

- **Courant :** C/C++, Python, Go, TensorFlow, Keras, LᴬTEX, VS Code, GNU/Linux, OS X, Git, SQL, MySQL.
- **Niveau avancé :** Raspberry Pi, Matlab, Wolfram Mathematica, JavaScript, PHP, HTML5, Java, Windows, MongoDB, OpenSSH, GnuPG, Apache Server, NumPy, GIMP.