

# Leopoldo Pla Sempere

SOFTWARE ARCHITECT · PART-TIME LECTURER

✉ b64encoded: bGVvcG9sZG9wbGFzZW1wZXJlQGdtYWlsLmNvbQ== | 🏠 lpla.github.io | 📧 lpla | 📄 leopoldoplasempere

## Experience

### Part-time lecturer

San Vicente del Raspeig

UNIVERSITY OF ALICANTE

Oct. 2019 - present

- Teaching in **Mathematics and Computer Science** degrees
- Programming algorithmic fundamentals and advanced techniques for software specification
- Use real-world examples and challenges from Google, Amazon, Microsoft and Spotify.

### Senior Data Scientist, Natural Language Processing and Machine Learning engineer

San Vicente del Raspeig

UNIVERSITY OF ALICANTE

Nov. 2017 - present

- Lead developer in **EU funded projects** in coordination with
  - University of Edinburgh, Prompsit, TAUS, Omnisien and John Hopkins University:
    - \* **ParaCrawl** (Provision of Web-Scale Parallel Corpora for Official European Languages)
    - \* **Paracrawl 2** (Broader Web-Scale Provision of Parallel Corpora for European Languages)
    - \* **Paracrawl 3** (Continued Provision of Web-Scale Parallel Corpora for European Languages) (paracrawl.eu)
  - University of Groningen, Jožef Stefan Institute, Prompsit:
    - \* **MaCoCu**: Massive collection and curation of monolingual and bilingual data: focus on under-resourced languages (macocu.eu)
- My work focus on **creating parallel corpora** for EU official languages **by a broad web crawling effort**.
- Coordinated the development of the **Open Source suite Bitextor** [bitextor/bitextor on Github] (bash, python, C++, snakemake, awk...), the production-ready tool that merges all this necessary technology chain for these CEF projects, run under **HPCC-Cloud scale** (Cambridge Service for Data-Driven Discovery, CSD3).
- 6 people team under my supervision.
- Cloud resources provided by **Amazon (AWS S3), Microsoft (Azure) and Docker**.
- OS administration of various Linux-based systems (Scientific Linux, Ubuntu, Fedora, Debian)

### Data Scientist, Natural Language Processing and Machine Learning engineer

Elche

PROMPSIT LANGUAGE ENGINEERING

May 2015 - Oct. 2017

- Developing natural language related projects as Reverso Context, processing full range of corpora and language resources
  - developed for Reverso (Softissimo, approx. rank 300 in Global Alexa Ranking)
  - **several million page-views per day**, big data Cloud scale computing
  - grown in an agile development team
  - using technologies as SVN, pandas, jupyter, bash (UNIX tools), word2vec, keras, scipy and Apache Solr
  - Language resources I was developing were in all available languages in the app (Arabic, German, English, French, Hebrew, Italian, Japanese, Dutch, Polish, Portuguese, Romanian and Russian)
  - High responsibilities in product releases and Quality Assessment.

### Full-stack developer intern

Elche

PROMPSIT LANGUAGE ENGINEERING

Feb. 2013 - May 2013

- Developing an internal news sentiment analysis suite
- Python and bash scripting for news crawling and plain text extraction from all downloaded HTML data
- AJAX, jQuery, CSSv3 and PHP for real-time interactive interface for dataset generation
- MySQL for dataset storage

## Education

### Master's Degree in Artificial Intelligence, Pattern Recognition and Digital Imaging

Valencia

POLYTECHNIC UNIVERSITY OF VALENCIA

Sep. 2014 - Jun. 2015

- GPA: 8.81/10
- M.Sc. specialised in **natural language processing and deep learning**
- Coursed machine learning MOOCs of Andrew Ng (Stanford University) and Yaser Abu-Mostafa (Caltech)
- Master Thesis: Audio classical composer identification in MIREX 2015: submission based on Structural Analysis of Music

### Computer Science Degree

San Vicente del Raspeig

UNIVERSITY OF ALICANTE

Sep. 2010 - Jun. 2014

- GPA: 8.53/10
- B.Sc. specialized in Informatics (robotics, machine learning, artificial vision, compilers)
- **Enrolled in the group of high profile academics programme**
- Diploma Thesis: Dodecaphonic music composer assistant with OpenMusic

## Publications & Awards

### Building Domain-specific Corpora from the Web: the Case of European Digital Service Infrastructures

LREC 2022

van Noord et al.

Marseille, France

June, 2022

### MaCoCu: Massive collection and curation of monolingual and bilingual data: focus on under-resourced languages

EAMT 2022

Bañón et al.

Ghent, Belgium

June, 2022

### ParaCrawl: Web-Scale Acquisition of Parallel Corpora

ACL 2020

Bañón et al.

Online

6 July, 2020

### Computer Science Degree graduation and M.Sc scholarship

BEST ACADEMIC RECORD

University of Alicante, Generalitat Valenciana

San Vicente del Raspeig

2015

### B.Sc thesis defense

HIGHEST MARK (OUTSTANDING WITH HONORS)

University of Alicante

San Vicente del Raspeig

2014

## Presentation & Associations

### EAMT

WEBMASTER FOR THE EUROPEAN ASSOCIATION OF MACHINE TRANSLATION

- Technical role managing services as eamt.org, lists.eamt.org, mt-archive.net and several EAMT congress websites

Remote

July 2021 - Present

### DockerCon 2022

PRESENTER FOR <THE BITEXTOR PIPELINE>

- Introduced the Bitextor pipeline and explained its uses with Docker

Remote

May, 2022

### Docker Community All-Hands #4 2021

PRESENTER FOR <THE BITEXTOR PIPELINE>

- Introduced the Bitextor pipeline and explained its uses with Docker, including a live demo

Remote

Dec. 2021

## Skills

### Natural Language Processing tools

Moses, Marian, Bleualign, Bitextor, Bicleaner

### Scripting and tools

git, CMake, tmux, Python, Lisp, Bash, AWK, Perl, systemd

### Machine Learning Toolkits

TensorFlow, Keras, PyTorch, pandas, sklearn, numpy, scipy, matplotlib

### Programming

C++, C, Java, Golang, JAVA, Rust, LaTeX, CMake, Snakemake

### Languages

English (proficient), Spanish (native), Catalan (native)

## Extracurricular Activity

### Clarinet Professional degree

PROFESSIONAL CONSERVATORY OF MUSIC OF ELCHE

- Specialized in jazz improvisation with Miguel García Ferrer
- MOOCs on jazz improvisation with Gary Burton (Berklee School of Music)
- MOOC on fundamentals of rehearsing music ensembles with Dr. Evan Feldman (University of North Carolina at Chapel Hill)
- After studies, gained experience in several youth orchestras, symphonic wind orchestras and rock, pop, swing and jazz bands, writing arrangements, improvising, while playing clarinet and piano as soloist. Non-professional experience with saxophone, guitar, ukulele, bass and drums.

Elche

Sep. 2004 - Jun. 2010

### Hardware-related projects

SELF-TAUGHT

- Open source PCB layout design using EAGLE and KiCad, with real production in factories as OSHPark and JLCPCB, for 80's computers.
- Electronic repairs and programming on vintage video-game platforms. Winner of the "best sound" award at "Retroconsolas Alicante 2013" for the Amstrad CPC 464.
- Open source 3D design (Autodesk), modeling and printing. Volunteer with "3D Makers Elche" group during COVID-19 lockdown by 3D-printing more than two hundred face-masks for near hospitals.