

Research engineer with 6+ years of experience working in academia on data science and engineering subjects seeking to employ the acquired skills for more practical applications. Well versed in the Python data libraries, thorough and independent-minded in the development and presentation of experimental results.

SKILLS

Tools and Languages	Python (Polars, Pandas, scikit-learn, Matplotlib), Git, $\LaTeX$ , shell
Fields of research	Tabular learning, data engineering, benchmarking, feature selection
Communication	English (professional level), Italian (native), French (conversational level)

TECHNICAL EXPERIENCE

<b>Research Engineer</b> <i>P16 – INRIA Saclay, Dataiku</i>	<b>Oct 2024 – current</b> <i>Saclay, France</i>
<ul style="list-style-type: none"><li>Development of Skrub, an open source tabular data processing library.</li></ul>	
<b>Postdoctoral researcher</b> <i>SODA Team – INRIA Saclay, Dataiku</i>	<b>Oct 2022 – September 2024</b> <i>Saclay, France</i>
<ul style="list-style-type: none"><li>GB-scale data cleaning and engineering.</li><li>Development of YADL, a synthetic benchmarking data lake.</li><li>Development of “Retrieve, Merge, Predict”, a pipeline for analyzing methods for automating table augmentation.</li></ul>	
<b>Doctoral Student</b> <i>EURECOM</i>	<b>Oct 2018 – April 2022</b> <i>Sophia Antipolis, France</i>
<ul style="list-style-type: none"><li>Development of EmbDI, a data integration algorithm based on embeddings of tabular data.</li><li>Development of GRIMP, a data imputation algorithm based on graph neural networks and multi-task learning.</li><li>Supervision of several master students as Teaching Assistant.</li></ul>	
<b>Software Developer</b> <i>SAP Labs</i>	<b>Mar 2018 – Aug 2018</b> <i>Mougins, France</i>
<b>Internship</b> <i>SAP Labs</i>	<b>July 2017 – Feb 2018</b> <i>Mougins, France</i>

EDUCATION

<b>PhD in Computer Science, Telecommunications and Electronics</b> , <i>Sorbonne Université, France</i>	2018–2022
<b>Master degree in Computer Security</b> , <i>EURECOM, France</i>	2016–2018
<b>Master degree in Communication and computer networks engineering</b> , <i>Politecnico di Torino, Italy</i> , 110/110	2015–2018
<b>Bachelor degree in Computer Engineering</b> , <i>Politecnico di Torino, Italy</i> , 107/110	2012–2015

ACTIVITIES

<b>Retrieve, Merge, Predict</b> , an experimental pipeline to benchmark table augmentation from data lakes.	2022–2024
<b>YADL (Yet Another Data Lake)</b> , a synthetic data lake for stress-testing SOTA augmentation methods.	2022–2024
<b>GRIMP (Graph embeddings for Relational data IMPutation)</b> , an imputation algorithm based on GNNs	2021–2022
<b>EmbDI (Embeddings for Data Integration)</b> , a system for generating table embeddings for data integration	2019–2022
<b>EDBT Summer School (Extracting Hidden Knowledge from Heterogeneous Massive Data)</b>	September 2019
<b>eROCK</b> , a clustering algorithm for categorical data	Fall 2017
<b>Semester project: studying the effect of Wifi networks on network protocols for IoT devices</b>	Spring 2017
<b>Marco Poli</b>	Spring 2015

PUBLICATIONS

Cappuzzo, R., Papotti, P., & Thirumuruganathan, S. (2020, June). **Creating embeddings of heterogeneous relational datasets for data integration tasks** – 2020 ACM SIGMOD.

Cappuzzo, R., Papotti, P., & Thirumuruganathan, S. (2021, September). **EmbDI: Generating Embeddings for Relational Data Integration** – 2021 SEBD.

Cappuzzo, R., Papotti, P., & Thirumuruganathan, S. (2024, March). **Relational Data Imputation with Graph Neural Networks** – 2024 EDBT.

Cappuzzo, R., Varoquaux, G., Coelho, A., & Papotti, P. (2024, May). **Retrieve, Merge, Predict: Augmenting Tables with Data Lakes** – Arxiv preprint.