

# Washington Cunha | AI Researcher

✉ waasluiz@gmail.com

🎓 <https://scholar.google.com.br/citations?user=TiRmr48AAAAJ>

🐙 [www.github.com/waashk](https://www.github.com/waashk)

🌐 [www.linkedin.com/in/washington-l-m-cunha/](https://www.linkedin.com/in/washington-l-m-cunha/)

- 📌 **Main Research:** Information Retrieval, Search and Ranking, Artificial Intelligence, and Machine Learning, focusing on Recommender Systems, Text Classification, and Feature Engineering.
- 📌 I am an Information Retrieval (IR) Researcher with experience in large language models (LLMs) and generative AI, specializing in transformer-based architectures. I have over eight years of experience in AI, proven through my work in academic research and industry. I recently spent a sabbatical as a visiting researcher at the CNR-ISTI in Italy under the supervision of Fabrizio Sebastiani.
- 📌 Currently, I'm focusing on reducing the training cost of LLMs, specifically generative models, through Instance Selection (IS). In sum, the IS goal is to reduce the training set size while maintaining the effectiveness of the trained models and reducing the training cost.
- 📌 International Collaborations: Berthier Ribeiro-Neto (Google); Fabrizio Sebastiani (ISTI-CNR); Nicola Ferro (UniPD); Andrea Esuli (ISTI-CNR); Davide Bacciu (UniPi); Antonio Carta (UniPi).

## Education

- 2019 – 2024 📌 **Ph.D. in Computer Science** at the Federal University of Minas Gerais.  
*A Comprehensive Exploitation of Instance Selection Methods for Automatic Text Classification*
- 2023 – 2024 📌 **Visiting Researcher** at the Consiglio Nazionale delle Ricerche, AI4Text, Italy.  
Project: *When Instance Selection meets Feature Selection* - Supervisor: Fabrizio Sebastiani
- 2018 – 2019 📌 **Master Degree in Computer Science** at the Federal University of Minas Gerais.  
Thesis title: *Extended pre-processing pipeline for text classification: On the role of meta-feature representations, sparsification and selective sampling*. **Honorable Mention** in the Masters Theses Contest of the Brazilian Database Symposium. **CTDBD – SBBD'21**
- 2014 – 2018 📌 **Bachelor Degree in Computer Science** at the Federal University of São João del-Rei.  
Thesis title: *A Feature-oriented Sentiment Rating for Mobile App Reviews*.

## Technical and Scientific Experience

### Head AI Researcher at AnaHealth (Full time: 03/2024 – Current)

- 📌 Leading the R&D team on proposing LLMs for Digital Primary Health Care to support AnaHealth's experts in personalizing care contextualizing based on individual patient records.
- 📌 Main achievements: The proposed model offers more cohesive and adjustable care, providing key information 5.2 times faster than reading the medical record, and respects each patient's individuality.

### AI Researcher at the Istituto di Scienza e Tecnologie dell'Informazione (09/2023 – 03/2024)

- 📌 Research, Development, and Application of instance selection approaches into the Generative AI data pipeline. (Automatic Text Classification and Quantification project)
- 📌 Main achievements: Published paper on ICTIR'24 on quantum computing; Under-review paper at the ACM TOIS (minor review); Ongoing research information retrieval with instance selection support.

### Data Science Consultant – CGEE - Brazilian government (Fixed-term contract - 2023)

- 📌 Monitoring Brazilian scientific and technological production to identify trending and emerging themes. Construction of indicators for key R&D variables to redistribute grants among researchers.
- 📌 Main achievements: This approach allowed the Brazilian government to optimize the allocation of resources, ensuring a higher return on investment by maximizing the expected research publications.

## Technical and Scientific Experience (continued)

### AI Researcher at ASTREIN/UFGM (Part time: 2021 – 2022)

- R&D of advanced and recent natural language processing strategies aimed to improve the tasks of association with description patterns, description retrieval, taxonomy classification, and association with external catalogs.
- Building an approach to effectively increase the efficiency and effectiveness of the company's end processes through the research, development, testing, adaptation, and application to retrieve complete material purchase specifications.
- Main achievements: Our solution was able to retrieve the correct standardized product in the top-5 ranking position in 71% of the cases and its correct category in the 1st position in 80% of the situations, providing a solution 3.7 times faster than the Astrein's previous approach.

### Data Scientist at Delloite (Full time: 2020 – 2021)

- Creating sophisticated analytics products by employing machine learning and AI.
- Main achievements: the primary outcome was the development of an innovative accident risk prediction model, Safety Analytics, implemented at Gerdau. This led to Gerdau receiving the Excellence in Health and Safety recognition from worldsteel, the leading global organization in the steel industry. The award acknowledged the company's success in significantly reducing the number of severe accidents, further demonstrated by achieving the lowest accident severity rate in the company's history.

## Awards and Achievements

- 2024 ■ **Best Reviewer Award - SIGIR'24**
  - **Ranked among the top ten Brazilian Scientific Initiation research projects** selected by the Brazilian Computer Society (as a co-advisor). **CTIC'24 – SBC.**
- 2023 ■ **CAPES-PRINT** for spending a sabbatical semester abroad.
  - **SIGIR Student Awards** for present an accepted full research paper at the SIGIR'23.
  - **Honorable Mention** in WFA – Tools and Applications Workshop – **WebMedia'23.**
- 2021 ■ **Honorable Mention** in the Masters Theses Contest of the Brazilian Database Symposium. **CTDBD – SBB'D'21**
- 2019 ■ **Ranked among the top ten Brazilian Scientific Initiation research projects** selected by the Brazilian Computer Society (as a co-advisor). **CTIC'19 – SBC.**

## Top six selected research publications

The full list can be seen in <https://scholar.google.com.br/citations?user=TiRmr48AAAAJ>.

- 1 Exploiting Contextual Embeddings in Hierarchical Topic Modeling and Investigating the Limits of the Current Evaluation Metrics. (2024). **Computational Linguistics – Imp. Fac.: 9.3.**
- 2 On representation learning-based methods for effective, efficient, and scalable code retrieval. (2024). **Neurocomputing – h-index: 196, Imp. Fac.: 5.5.**
- 3 A comparative survey of instance selection methods applied to nonneural and transformer-based text classification. (2023). **ACM Computing Surveys – Imp. Fac.: 23.8.**
- 4 A Quantum Annealing Instance Selection Approach for Efficient and Effective Transformer Fine-Tuning. (2023), In *The 14th international conference on the theory of information retrieval -ICTIR'24.*
- 5 An Effective, Efficient, and Scalable Confidence-Based Instance Selection Framework for Transformer-Based Text Classification. (2023), In **SIGIR'23 – h5-index: 75.**
- 6 On the class separability of contextual embeddings representations—or “the classifier does not matter when the (text) representation is so good!” (2023). **IP&M – Imp. Fac.: 7.4.**