

Nuno Gonçalves

✉ nuno.m.goncalves@tecnico.ulisboa.pt ☎ +351 963784966 🌐 nunomtg.github.io 📄 [nunomtg](#)

Education

Instituto Superior Técnico <i>Master of Science – Data Science and Engineering</i> <ul style="list-style-type: none">◦ GPA: 18.2/20.0	<i>Sept 2022 – July 2025</i> (Est.)
Technical University of Munich <i>Erasmus Program – Mathematics Department</i> <ul style="list-style-type: none">◦ GPA: 18/20	<i>Oct 2023 – Feb 2024</i>
Instituto Superior Técnico <i>Bachelor of Science – Physics</i>	<i>Sept 2019 – June 2022</i>

Experience

Deep Learning Research Intern <i>SARDINE Lab</i> <ul style="list-style-type: none">◦ Optimized the Entmax activation function for GPUs using Triton, achieving a 15x speed improvement.◦ Designed and implemented a custom algorithm inspired by FlashAttention2, enabling Entmax usage in attention layers for long-context Transformers.	<i>Lisbon, PT</i> <i>March 2024 –</i>
ML Engineer Intern <i>Talka.ai</i> <ul style="list-style-type: none">◦ Annotated data to incorporate new gestures into a vision-language model, enhancing its ability to classify gestures, behaviors, and speech for insights in sales meetings.◦ Designed and implemented a custom data augmentation pipeline, increasing recall by up to 26% for under-represented classes.	<i>Lisbon, PT</i> <i>Oct 2023 – Jan 2024</i>
Software Engineer Intern <i>CERN/LIP</i> <ul style="list-style-type: none">◦ Optimized the GPU-based version of a critical algorithm for identifying particle showers within the Large Hadron Collider (LHC) at CERN.◦ Improved the number of events processed per second by approximately 2x compared to the previous implementation.◦ Developed a custom GPU-agnostic heuristic to optimize kernel launch parameter selection for performance.	<i>Lisbon, PT</i> <i>Jan 2022 – July 2022</i>
Research Intern <i>INESC (Systems and Computing Engineering Institute)</i> <ul style="list-style-type: none">◦ Developed a Neural Network for signal classification capable of achieving 90% accuracy in detecting magnetic nanoparticles within a magnetic flow cytometry lab-on-a-chip device.	<i>Lisbon, PT</i> <i>Dec 2021 – March 2022</i>

Publications

Nuno Gonçalves, Marcos V. Treviso, André F. T. Martins. *AdaSplash: Adaptive Sparse Flash Attention*. To appear in *Proceedings of the 42nd International Conference on Machine Learning (ICML)*, 2025. [Spotlight](#).

Awards

IST Academic Excellence Award Awarded for achieving the highest academic distinction across all programs at Instituto Superior Técnico.	2024
Celfocus Academic Merit in Machine Learning Recognized as one of the top five students in the Machine Learning course for exceptional academic performance.	2023

Projects

Portuguese Word Embeddings



- Implemented and trained word models (HAL, CBOW, Skip-Gram) in a Portuguese corpus, creating contextualized vector representations of portuguese words.

Recommendation Systems In Pure NumPy



- Built and compared different recommendation systems from scratch, including the famous Funk's SVD and multiple non-negative matrix factorization (NMF) methods, fully implemented in pure NumPy.
- Predicted user ratings on the MovieLens dataset and visualized latent feature clustering.
- Explored connections between matrix factorization techniques and auto-encoders for latent representation.

Technologies

Languages: Python, Triton, C, C++, Rust, SQL

APIs and Technologies: CUDA, PyTorch, Git