

JOAN VELJA

+39 3319261802 joan.velja@student.uva.nl linkedin.com/joanvelja github.com/joanvelja
joanvelja.vercel.app scholar.google.com/joanvelja

EDUCATION

University of Amsterdam

Exp. July 2025

M.S. Artificial Intelligence (Current GPA: 8.6/10, equivalent to 4.0/4.0, on track for Cum Laude graduation) Amsterdam, Netherlands

- **Relevant Coursework:** Machine Learning I, Natural Language Processing I, Deep Learning I-II, Foundation Models
- **Current research interests:** AI Alignment, Reinforcement Learning

University of Technology Sydney (Exchange Semester)

June 2023

B.Sc Artificial Intelligence (GPA: 4.0/4.0)

Sydney, Australia

- Selected by academic merit to participate in the Exchange Program 2022-2023 with a full ride scholarship.
- **Relevant Coursework:** Deep Learning and Convolutional Neural Network, Natural Language Processing, Cloud computing and software as a service, Data Structures and Algorithms

Università Bocconi

July 2023

B.Sc Data Science (GPA: 29.3/30, equivalent to 4.0/4.0); Graduated 110/110 Cum Laude

Milan, Italy

- **Relevant Coursework:** Fundamentals of Computer Science, Computer Programming, Machine Learning, Mathematics, Advanced Mathematics, Statistics, Advanced Statistics, Big Data and Databases, Game Theory
- **Awards:** Awarded for 3 years scholarship for deserving students, top 2% in my cohort.
- **Thesis:** An unorthodox shift in the variance-bias tradeoff in Neural Networks: the double descent phenomenon and the ease of training in the overparametrized regime, grade: 4/4. Read about it [here](#).
- Thesis Supervisor: **Enrico Maria Malatesta**

EXPERIENCE

MSc Honors Programme - University of Oxford

Exp. Jan 2025 – Jun 2025

Thesis

Oxford, United Kingdom

- * Visiting researcher at University of Oxford.
- * Thesis topic: Instructing Strong Provers with Weak Verifiers
- * Investigating alternative reward functions to mitigate the trade-off between model accuracy and human-checkability in prover-verifier games. Exploring smoother reward formulations and techniques like curriculum learning to improve correctness and legibility, reducing the alignment tax.
- * Examining whether non-zero-sum training dynamics allow implicit collusion between provers and verifiers, leading to reliance on spurious correlations instead of genuine legibility. Designing experiments to assess this vulnerability and evaluating whether a zero-sum setup could improve adversarial robustness and legibility.
- * Introducing sequential play by allowing deceptive provers to condition responses on honest prover outputs. Investigating whether this setup strengthens adversarial challenges and forces provers to generate more robustly reasoned and legible solutions, as well as improving verifier effectiveness in distinguishing between genuine and misleading proofs.
- * Supervisor: **Alessandro Abate**.

LASR Labs

Jul 2024 – Oct 2024

AI Alignment Researcher

London, United Kingdom

- * 12 weeks program aimed at producing a paper and accompanying blog post that makes incremental progress on an important problem in AI safety.
- * Project revolved around elicitation of *steganography* in multi-agent AI systems under optimization pressure (In Context Learning, Reinforcement Learning) drawing from the AI Control research agenda proposed by Redwood Research and Model Organisms by Anthropic.
- * Goal was to assess whether steganography could emerge as an instrumental goal and to red-team paraphrasing, exploring alternative mitigation strategies.
- * First paper to demonstrate instrumental emergence of steganography.
- * Supervisors: **Nandi Schoots, Dylan Cope, Charlie Griffin**, and several feedback rounds with **Fabien Roger** (Anthropic)
- * Find our preprint, submitted to ICML 2025, [clicking here](#) and slides to presentation [clicking here](#).

1000 Kelvin GmbH

R&D Software Engineer

Aug 2023 – Dec 2023

Berlin, Germany

- * Developed a Large Language Model (Llama-2 Finetune) specialized in Physics and Additive Manufacturing.
- * Improved the company's proprietary model training pipeline, through a new approach towards generation and utilization of synthetic data.

World Food Programme

Data Scientist - Intern

Jun 2022 – Sep 2022

Rome, Italy

- * Developed analytical dashboards in Tableau, still in use in the Head of Business Continuity department, by leveraging unused organizational data.
- * Applied predictive classification algorithms to dashboard data, resulting in improved accuracy in predicting country adequacy levels, to enhance organizational strategic planning.
- * Conducted a comprehensive comparative analysis of countries in North West Africa from a Business Continuity perspective, leading to identification of key operational gaps or best practices, which informed the fine-tuning of dashboard use cases and was presented to the heads of the Regional Bureau in Dakar, Senegal.

PROJECTS

Pitfalls from Partial Observability in RLHF | *Reinforcement Learning, Stable Baselines 3*

- * In recent work, a theoretical analysis of the consequences of partial observability in RLHF has been made: agents may deceptively inflate performance, or overjustify behavior to make an impression. In this project, I repeated the analysis of the paper by Lang et. al (2024), but focusing on empirical demonstrations of the failure modes and solutions, thus complementing the existing theoretical analysis. I've done so by implementing a synthetic human in toy tasks to make "preference choices" based on partial observations, by e.g. modeling the human as Boltzmann rational. Moreover, I conducted variations on partial observability frameworks, with different experimental setups.
- * Find the output of this project [here](#).
- * Supervision by: **Leon Lang** (University of Amsterdam), **Davis Foote** (Center for Human-Compatible AI, UC Berkeley).

Dynamic Vocabulary Pruning in Early-Exiting LLMs | *Large Language Models, Pytorch*

- * This project explores the effectiveness of early exiting strategies in Large Language Models (LLMs) during inference to reduce computational overhead. By allowing the model to exit at specific points in the attention stack based on confidence thresholds, we seek to maintain performance while reducing runtime. The core idea of the project is to speed up the *unembedding* operation during inference by projecting onto a subset of dimensions. This is of interest as a natural advancement to the original work of Schuster et al. 2022. The project, accepted as a workshop paper at ENLSP-IV NeurIPS 2024, shows how we can attain significant gains in FLOPs and latency, while allowing the model to exit several significantly earlier and retaining up to 100% of its performance.
- * Find the paper – submitted to ENLSP-IV NeurIPS 2024 – [here](#).
- * Supervision by: **Metod Jazbec**.

PUBLICATIONS

Golechha S.*, Chaudhary M.*, **Velja Joan***, Abate A., Schoots N. 'Modular Training of Neural Networks aids Interpretability', *Under Review*

Mathew Y.*, Matthews O.*, McCarthy R.*, **Velja Joan***, Schroeder de Witt C., Cope D., Schoots N. 'Hidden in Plain Text: Emergence & Mitigation of Steganographic Collusion in LLMs', *Under Review*

Mathew Y.*, Matthews O.*, McCarthy R.*, **Velja Joan***, Cope D., Schoots N. 'Emergence of Steganographic Collusion between Large Language Models', *Towards Safe & Trustworthy Agents, NeurIPS 2024 Workshop*

Velja Joan*, Abdel Sadek K.*, Nulli M.*, Vincenti J.*, Jazbec M. *Dynamic Vocabulary Pruning in Early-Exit LLMs*, *ENLSP-IV NeurIPS Workshop 2024*

Velja Joan*, Abdel Sadek K.*, Nulli M.*, Vincenti J.* 'Explaining RL Decisions with Trajectories: A Reproducibility Study', *TMLR 2024*

Velja Joan*, Divak Adam* 'Assessing expertise overlap in Mixture of Experts Architectures', MechInterp Hackathon Project

Velja Joan, 'On the importance of interpretable code - Einops', Blogpost

AWARDS, HONORS AND SCHOLARSHIPS

Fellowship Arcadia Impact, LASR Labs: AI Safety Research Intern Summer '24, 15,000 \$	2024
Scholarship and Travel Grant Human Aligned AI Summer School, 1,000 \$	2024
Scholarship Università Bocconi: Merit Scholarship, 20,000€	2020, 2021, 2022
Award Brembo Sensify Hackathon, 3rd Place, 2,000€	2022
Scholarship Comune di Schio: Merit Scholarship, top 5% of cohort for the class of 2019, 1,000€	2019
Scholarship Banca Intesa San Paolo: Merit Scholarship for outstanding academic results, 500€	2015, 2016, 2017, 2018
Scholarship Regione Veneto: Merit Scholarship for outstanding academic results, 300€	2014

EXTRACURRICULARS

Co-founder Bocconi AI & Neuroscience Student Association	2022
Team Formation, Lead Manager Think Tank Tortuga	2021-23
Team Captain Rugby Alto Vicentino	2016-19

TECHNICAL SKILLS

Programming Languages: Python, C++, R, Julia (proficient) - C#, HTML, SQL, Stata, Tableau (advanced) - CSS (familiar)
Frameworks: Vue.js, Django, TensorFlow, PyTorch, TransformerLens, Gym, StableBaselines, d3rlpy, Transformers, TRL
Languages: Albanian, Italian (Native), English (Fluent), German (Intermediate)

LINKS

Bachelor's Thesis: https://joanvelja.vercel.app/static/pdf/BSc_Thesis.pdf
Hidden in Plain Text (Steganography paper): <https://arxiv.org/abs/2410.03768>
PORLHF : <https://joanvelja.vercel.app/static/pdf/PORLHF.pdf>
Dynamic Pruning : <https://arxiv.org/abs/2410.18952>