# Bálint Gyevnár

Edinburgh, UK | balint.gyevnar@ed.ac.uk | gbalint.me

## RESEARCH INTERESTS

PhD student focusing on human subjects testing of AI safety and explainable multi-agent autonomous systems for trustworthy human-agent collaboration, with applications to autonomous vehicle planning.

## SKILLS

**Programming:** Python (PyTorch, Transformers, Pandas, etc.), R (dplyr, ggplot2, rlmer, etc.), C++ (CARLA), C#, Bash, Java, Haskell;

**Data analysis:** human subjects studies, unsupervised topic modeling, mixed effects regression, statistical hypothesis testing, data visualization;

**Languages:** English (fluent), German (fluent), Japanese (intermediate), Hungarian (native).

## EDUCATION

**University of Edinburgh**　　　　　　　　　　　　　　　　　　Sep. 2021 – May 2025 (est.)
*PhD in Natural Language Processing with Integrated Studies*　　　　　　　　*Edinburgh, UK*
*Supervisors: Stefano Albrecht, Shay Cohen, Christopher Lucas*

**University of Edinburgh**　　　　　　　　　　　　　　　　　　　Sep. 2016 – May 2021
*Integrated Master of Informatics*　　　　　　　　　　　　　　　　　　*Edinburgh, UK*
*Supervisor: Maria Wolters*

## PROJECTS

**Bridging shared research challenges amid responsible AI wars**　　　　Jul. 2024 – present
- Curation of corpus of 3K+ papers on AI safety and AI ethics;
- Qualitative data analysis and visualization (e.g., coding, graph analysis);
- Quantitative unsupervised topic modeling and analysis (e.g., BERTopic);
- Co-authoring with Shannon Vallor and Atoosa Kasirzadeh.

**Causal Explanations for Sequential Decision-Making in Multi-Agent Systems**　　　　Sep. 2021 – present
- Counterfactual reasoning with RL planning for causally-grounded explanations in natural language;
- Two large-scale human subjects studies to evaluate natural and automatically generated explanations;
- Curation of HEADD: The Human Explanations for Autonomous Driving Decisions dataset;
- Integration of LLMs with CEMA in an RAG approach to improve the quality of explanations.

## EXPERIENCE

**Research Assistant**　　　　　　　　　　　　　　　　　　　　Jul. 2023 – present
*University of Edinburgh*　　　　　　　　　　　　　　　　　　　　*Edinburgh, UK*
- Researching the intersection of AI safety and AI ethics to build bridges among research problems;
- Large scale quantitative literature analysis with unsupervised natural language processing tools;
- Curation, topic coding, and qualitative analysis of large corpora of papers.

**Research Internship**　　　　　　　　　　　　　　　　　　　　May 2020 – Oct. 2020
*Five AI Ltd.*　　　　　　　　　　　　　　　　　　　　　　　　*Edinburgh, UK*
- Development and evaluation of goal-based interpretable prediction and planning for autonomous vehicles;
- Scenario-based and open-world testing and results collection;
- Main contributor of open-source implementation on GitHub with added support for CARLA.

**Teaching Assistant**                                                    Sep. 2019 – present
*University of Edinburgh*                                                    *Edinburgh, UK*

- Teaching assistant for "Evaluating Sustainable Lands & Cities" and "Data Mobility & Infrastructure";
- Supervision of master's students and tutor for ~12 students for machine learning;
- Marker for courses in natural language processing, reinforcement learning, and machine learning.

## VOLUNTEERING

**Sports Club Executive Member**                                          Sep. 2022 – present
*Edinburgh University Volleyball Club*                                       *Edinburgh, UK*

- (2024-25; Secretary) Public outreach and networking with alumni members and organizing an event series;
- (2023-24; VP) Large-scale events, public speaking, timetabling, HR management of 220+ members;
- (2022-23; Treasurer) Setting up an annual budget, and managing a cash flow of £70k.

## AWARDS

**Colours Award for Outstanding Volunteering Contribution to Sports**         Jun. 2024
*Edinburgh University Sports Union*                                          *Edinburgh, UK*

**AI100 Early Career Essay Competition Featured Essay**                       Aug. 2023
*One Hundred Year Study on Artificial Intelligence (AI100)*               *Stanford University*

**Trustworthy Autonomous Systems Early Career Researcher Award**              Jun. 2023
*4,000 GBP; UK Research & Innovation*                                       *Southampton, UK*

**Shape the Future of ITS Competition; 3rd Place**                            Aug. 2022
*1,000 USD; IEEE Intelligent Transportation Systems Society*                          *USA*

## SELECTED PUBLICATIONS

**Conference:**
- **People Attribute Purpose to Autonomous Vehicles When Explaining Their Behavior**
  *[under review] ACM Conference on Human Factors in Computing Systems, CHI 2025;*
  B. Gyevnar, S. Droop, T. Quillien, S.B. Cohen, N.R. Bramley, C.G. Lucas, S.V. Albrecht.
- **Causal Explanations for Sequential Decision-Making in Multi-Agent Systems**
  *23rd International Conference on Autonomous Agents and Multiagent Systems, AAMAS 2024;*
  B. Gyevnar, C. Wang, S.B. Cohen, C.G. Lucas, S.V. Albrecht.
- **Bridging the Transparency Gap: What Can Explainable AI Learn From the AI Act?**
  *26th European Conference on Artificial Intelligence, ECAI 2023;*
  B. Gyevnar, N. Ferguson, B. Schafer.
- **Interpretable Goal-based Prediction and Planning for Autonomous Driving**
  *International Conference on Robotics and Automation, ICRA 2021;*
  S.V. Albrecht, C. Brewitt, J. Wilhelm, B. Gyevnar, F. Eiras, M. Dobre, S. Ramamoorthy.
- **GRIT: Fast, Interpretable, and Verifiable Goal Recognition with Learned Decision Trees for Autonomous Driving**
  *IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS, 2021;*
  C. Brewitt, B. Gyevnar, S. Garcin., S.V. Albrecht.

**Journal:**
- **Explainable AI for Safe and Trustworthy Autonomous Driving: A Systematic Review**
  *IEEE Transactions on Intelligent Transportation Systems, 25 (12), 19342-19364, IEEE T-ITS 2024;*
  A. Kuznietsov*, B. Gyevnar*, C. Wang, S. Peters, S.V. Albrecht. [* equal contribution]

**Workshop:**
- **A Human-Centric Method for Generating Causal Explanations in Natural Language for Autonomous Vehicle Motion Planning** [best paper runner-up]
  *Workshop on Artificial Intelligence for Autonomous Driving, IJCAI 2022;*
  B. Gyevnar, M. Tamborski, C. Wang, C.G. Lucas, S.B. Cohen, S.V. Albrecht.