

Dayyán O'Brien

dayyanobrien.github.io

+44 7770 655794

dayyanobrien@icloud.com

EDUCATION

University of Edinburgh	Edinburgh, UK
<i>PhD in Designing Responsible NLP</i>	<i>Sep 2025 – Aug 2029</i>
<ul style="list-style-type: none">• Supervision: Dr. Emily Allaway• Fully funded by the G-Research PhD Scholarship for outstanding academic merit.	
<i>Master of Informatics [transcript] (First Class Honours)</i>	<i>Sep 2019 – May 2024</i>
<ul style="list-style-type: none">• Dissertation supervision: Prof. Mirella Lapata• Courses: Natural Language Processing, Machine Learning, Reinforcement Learning, Calculus, Linear Algebra	
Abbey Christian Brothers' Grammar School	Newry, UK
<i>Secondary school (AAA in A-Level Mathematics, Physics and Computer Science)</i>	<i>Sep 2012 – May 2019</i>

PUBLICATIONS

HPLT 3.0: Very Large-Scale Multilingual Resources for LLM and MT. Mono- and Bi-lingual Data, Multilingual Evaluation, and Pre-Trained Models

Stephan Oepen, Nikolay Arefev, et al. (incl. Dayyán O'Brien)

arXiv:2511.01066. Under review for the Language Resources and Evaluation Conference (LREC). 2025.

MatheMagic: Generating Dynamic Mathematics Benchmarks Robust to Memorization

Dayyán O'Brien, Pinzhen Chen, Barry Haddow, Emily Allaway

arXiv:2510.05962. Under review for ACL Rolling Review (ARR). 2025.

MathNLP Workshop, Empirical Methods in Natural Language Processing (EMNLP) (non-archival). 2025.

DocHPLT: A Massively Multilingual Document-Level Translation Dataset

Dayyán O'Brien, Bhavitya Malik, Ona De Gibert Bonet, Pinzhen Chen, Barry Haddow, Jörg Tiedemann

Proceedings of the Conference on Machine Translation (WMT). 2025.

An Expanded Massive Multilingual Dataset for High-Performance Language Technologies

Laurie Burchell, Ona de Gibert, et al. (incl. Dayyán O'Brien)

Proceedings of the Association for Computational Linguistics (ACL). 2025.

Mind the Gap: Diverse NMT Models for Resource-Constrained Environments

Ona De Gibert Bonet, Dayyán O'Brien, Dušan Variš, Jörg Tiedemann

Proceedings of the Nordic Conference on Computational Linguistics (NoDaLiDa). 2025.

EMMA-500: Enhancing Massively Multilingual Adaptation of Large Language Models

Shaoxiong Ji, Zihao Li, et al. (incl. Dayyán O'Brien)

arXiv:2409.17892. Under review for the Journal of Data-centric Machine Learning Research (DMLR). 2024.

Prompting Numerical Commonsense Reasoning across Languages

Dayyán O'Brien

Outstanding Honours Thesis, School of Informatics, University of Edinburgh. 2024.

Numerical Commonsense Reasoning across Languages

Dayyán O'Brien

Outstanding Honours Thesis, School of Informatics, University of Edinburgh. 2023.

RESEARCH EXPERIENCE

University of Edinburgh	Edinburgh, UK
<i>Doctoral Researcher, supervised by Emily Allaway</i>	<i>Sep 2025 – Present</i>
<ul style="list-style-type: none">• Improving the compositionality of language models.	
<i>Research Assistant, supervised by Barry Haddow</i>	<i>Aug 2024 – Aug 2025</i>
<ul style="list-style-type: none">• Co-developed the HPLT 2.0 bitexting pipeline to mine parallel data.	
<i>Junior Research Assistant, supervised by Pinzhen Chen</i>	<i>May 2024 – Jul 2024</i>
<ul style="list-style-type: none">• Cleaned and processed data for EMMA-500, a 7B parameter LLM for over 500 languages.	
<i>Junior Research Assistant, supervised by Mirella Lapata</i>	<i>Jun 2022 – Apr 2024</i>
<ul style="list-style-type: none">• Curated and evaluated mNumersense: 36k+ Arabic, Chinese, Russian numeric commonsense sentences.	

AWARDS & ACHIEVEMENTS

- G-Research PhD Scholarship (2025 – 2029)**
- Outstanding Honours Project (2023 & 2024)**
- Runner-up Best Coursework for Reasoning and Agents (2021)**
- Exemplary Project for Foundations of Data Science (2021)**

MAJOR CONTRIBUTIONS & OPEN SOURCE

- Leader, Document-Level HPLT Corpus:** Large-scale, document-level parallel corpus (huggingface.co/datasets/HPLT/DocHPLT).
- Contributor, EMMA-500:** An open-source LLM for 500+ languages (github.com/MaLA-LM/emma-500).
- Contributor, HPLT 2.0 Bitexting Pipeline:** Parallel data mining pipeline (hplt-project.org/datasets/v2.0).

TEACHING EXPERIENCE

University of Edinburgh	Edinburgh, UK
<i>Demonstrator for Accelerated Natural Language Processing (INFR11125)</i>	Sep 2023 – Nov 2023
<i>Tutor for Foundations of Natural Language Processing (INFR10078)</i>	Jan 2023 – May 2023
<i>Tutor for Foundations of Data Science (INFR08030)</i>	Sep 2022 – May 2023
<i>Demonstrator for Foundations of Natural Language Processing (INFR10078)</i>	Jan 2024 – Mar 2024
Self-employed	Edinburgh, UK
<i>Private tutor</i>	Sep 2020 – Mar 2024

INDUSTRY EXPERIENCE

Kainos	Belfast, UK
<i>Software Engineer</i>	Mar 2018
Bombardier	Belfast, UK
<i>Engineer</i>	Sep 2017
Computer Hospital	Newry, UK
<i>Computer technician</i>	Oct 2017

LEADERSHIP & OUTREACH

G-Research Brand Ambassador	Edinburgh, UK
<i>G-Research</i>	Sep 2025 – Present
Committee Member, Community Events Planning	Edinburgh, UK
<i>Edinburgh Bahá'í Community</i>	Mar 2025 – Present
Primary School Micro:bit Programming Workshop Facilitator	Bathgate, UK
<i>University of Edinburgh</i>	Jan 2024 – Feb 2024

SKILLS

- Programming:** Python, Java, Haskell, L^AT_EX, SQL
- Libraries:** PyTorch, Tensorflow, HuggingFace, NLTK, NumPy, Slurm, Kubernetes, Docker, MTurk, pandas, sklearn, statsmodels, Matplotlib, Festival, HTK, Tkinter, Google Cloud VM, Weather & Maps API, Seaborn, Kivy, JUnit, Maven
- Languages:** English (Native), German (Professional)