

# Barkavi Sundararajan

email | LinkedIn | Google Scholar  
Aberdeen, United Kingdom

## PERSONAL SUMMARY

Final-year PhD researcher and AI data scientist specialising in NLP and LLMs, with hands-on experience building and evaluating end-to-end LLM pipelines for sensitive public-sector data. I develop scalable Python and SQL workflows on AWS (Athena, S3, Bedrock), enforce structured outputs with Pydantic validation, and design evaluation frameworks combining secure human validation and LLM as a judge to support reliable operational use and monitoring.

## TECHNICAL SKILLS

- **LLMs & Applied GenAI:** model selection, fine-tuning, prompt engineering, structured generation, LLM APIs, vector databases (ChromaDB), long-context handling (NBA play-by-play JSONs up to 80K tokens), evaluation-led iteration, agentic AI (PoC)
- **Safety & Guardrails:** PII redaction, schema enforcement (Pydantic), regex/spaCy post-processing, compliant LLM design
- **Evaluation & Validation:** human annotation protocols, LLM-as-Judge, error taxonomies, inter-annotator agreement (Cohen's/Fleiss'  $\kappa$ ), precision/recall, qualitative error analysis
- **Programming & Data:** Python (5+ years), SQL (Athena/PostgreSQL), Pandas, NumPy, Scikit-learn, PyTorch, Hugging Face, spaCy
- **Cloud & Delivery:** AWS (Bedrock, Athena, S3), Docker, Azure DevOps (CI/CD), web app deployment
- **Tools & Collaboration:** GitHub, BitBucket, JIRA, Confluence, Agile workflows

## WORK EXPERIENCE

### Ministry of Justice

Data Scientist Turing Intern

Nov 2024 – present

- Designed and implemented an end-to-end LLM extraction pipeline on the secure MoJ Analytical Platform; processed **150k** OASys free-text assessment notes and converted five free-text fields into 30 structured variables (e.g., domestic abuse flags, offence location types, anti-social behaviour, victim and weapon details) for **Safer Streets UK** analysis.
- Applied targeted sampling to focus on single-offence offenders with high/very-high risk scores (2022 to 2024), querying via SQL in **AWS Athena** (pydbtools) and cleaning and preprocessing meaningful records for LLM extraction.
- Built a modular Python pipeline calling AWS Bedrock LLMs, enforced structured outputs with **Pydantic** JSON schemas, and implemented **guardrails** using spaCy and regex for cleaning, normalisation, PII redaction and category grouping.
- Enabled secure human-in-the-loop validation by quickly upskilling in R and developing an **RShiny** validation app so social researchers could validate outputs directly on-platform without downloading sensitive data; recognised with a **MoJ Reward & Recognition** award for this work.
- Computed a **weighted accuracy metric** from **human validation** on 500 samples across 30 fields, achieving **75%–98%** weighted accuracy on **26/30** fields, then used **human-in-the-loop feedback** to refine the offence location extraction logic, improving weighted accuracy for offence location variables from **82.5% to 91.8%**.
- Built RShiny dashboards to visualise validation outcomes and summary analysis for social research stakeholders.
- Collaborated with Data First and BOLD stakeholders to align use cases and ensure the pipeline could scale to historical data processing.
- **Technologies:** Python, SQL, AWS Athena, S3, Bedrock, Pandas, Pydantic, spaCy, BERTopic, regex, GitHub, RShiny, Confluence

### University of Aberdeen

PhD Research student and Teaching Assistant

Oct 2021 – present

- Fine-tuned models including T5, FLAN-T5, BLOOM, and LLAMA 2-7B on tabular datasets (ToTTo, E2E, NBA Play-by-Play, and other custom tabular data). Curated zero-shot, few-shot, and Chain-of-Thought (CoT) prompts for LLAMA 3.1, Qwen 2.5, and Mistral-7B to handle diverse tabular input formats.
- Designed and demonstrated Practical classes for the courses: '**Natural Language Generation**', '**Data Mining with Deep Learning**' and '**Evaluation of AI Systems**'.
- Mentored five MSc AI students on their theses in Natural Language Processing, alongside my supervisor.
- Authored multiple peer-reviewed publications on factual accuracy and hallucination reduction in data-to-text generation (INLG 2025, NAACL 2024, CSL 2023, and GEM 2022).

## **James Fisher Asset Information Services**

*Full Stack Software Developer Intern*

*Nov 2020 – April 2021*

- Designed and developed an STS web application using **Django** and **MongoDB** for the backend, and JavaScript, Bootstrap, and CSS for the front-end to improve user experience. This app streamlined business data previously managed manually in static spreadsheets across multiple silos.
- Successful migration of the complete **Production Data** from static spreadsheets to the database-based application and optimised the back-end query performance.
- Migrated the media files from the code repository to Azure blob storage account with SAS token to ensure **secured access** to the blobs.
- Set up the **CI/CD pipelines** from scratch in **Azure DevOps** by following agile principles and deploying the web application in Microsoft Azure.
- **Technologies and Tools Used:** Python, Django, MongoDB, Azure DevOps, Microsoft Azure, BitBucket, GitHub

## **Plintron Mobility Solutions Pvt Ltd**

*IT and Strategic Accounts Consultant*

*Aug 2011 – Dec 2019*

- Led the **GDPR compliance project** for millions of telecom subscriber records, designing requirements for SAR processing, profiling restrictions, archival, encryption and reporting in collaboration with solution architects and product teams.
- Prepared **telecom churn** analyses and KPI-based service management packs for monthly reviews with European B2B clients, supporting data-led account decisions.
- **Tools:** Excel (data analysis), Salesforce, Clarizen, JIRA, Confluence; **Certifications:** AWS Business Professional, ITIL v3 Foundation.

## **PINC Wealth**

*Analyst Intern*

*Jan 2011 – May 2011*

- Analysed the correlation between BSE SENSEX and DOW JONES using the Markowitz model, Sharpe ratio, and correlation coefficient, and surveyed investor preferences. Applied **predictive analytics techniques** to provide insights into market dynamics and investor behaviour.

## **EDUCATION HISTORY**

---

<b>PhD in Computing Science</b>	Oct 2021 – present
<i>University of Aberdeen, United Kingdom</i>	
<b>MSc in Information Technology, With Distinction</b>	Jan 2020 – Jan 2021
<i>University of Aberdeen, United Kingdom</i>	
<b>Master of Business Administration, First class</b>	Sep 2009 – June 2011
<i>University of Madras, India</i>	
<b>Bachelor of Technology in Computer Science Engineering, First class</b>	Sep 2005 – May 2009
<i>Pondicherry University, India</i>	

## **RESEARCH WORK**

---

<b>PhD Research Student and Teaching Assistant</b>	October 2021 to present
• Presented the paper ' <b>Input Matters: Evaluating Input Structure's Impact on LLM Summaries of Sports Play-by-Play</b> ' at the INLG 2025 conference, Vietnam.	
• Presented the paper ' <b>Improving Factual Accuracy of Neural Table-to-Text Output by Addressing Input Problems in ToTTo</b> ' at the NAACL 2024 main conference, Mexico (a top-tier NLP venue with an acceptance rate of approximately 20%).	
• Collaborated with a Senior PhD Researcher on a journal paper titled ' <b>Evaluating factual accuracy in complex data-to-text</b> ' for basketball summaries, published in the Computer Speech & Language journal in 2023.	
• Presented the paper ' <b>Error Analysis of ToTTo Table-to-Text Neural NLG Models</b> ' at the Generation, Evaluation & Metrics (GEM) workshop in EMNLP 2022, Abu Dhabi.	

## **ACADEMIC ACTIVITIES AND ACHIEVEMENTS**

---

- Awarded **Scottish Informatics and Computer Science Alliance (SICSA)** PhD Travel Funding to support conference travel and the presentation of my research at INLG 2025.
- Awarded as a Virtual Volunteer for the **North American Chapter of the Association for Computational Linguistics (NAACL)** 2024 conference. Volunteered at multiple events organised by the **SICSA** and the **Alan Turing Institute**.
- Student Representative for MSc IT cohort (2020-2021): Represented student concerns to Computing Staff during the transition to blended learning in 2020.
- Presented 'Changes in Retail Banking' at the International Conference on Management of Change (ICMOC), 2010.
- Presented 'Pattern analysis using Neural Networks with Back Propagation algorithm for recognising fingerprints and handwritten digits' in National level Technical Symposium, 2009.
- Secured a CGPA of 10 out of 10 in Engineering Graphics in the Engineering University exam, 2005.