

Suchir Salhan

Machine Learning PhD at the University of Cambridge & Starred First Cambridge Graduate

✉ sas245@cam.ac.uk 🌐 www.suchirsalhan.com 📧 @suchirsalhan 🏠 Suchir Salhan 📍 Cambridge, UK

First Year PhD Candidate pre-training performant Small Language Models for Interactive Language Learning. I work with state-space models and hybrid Transformer architectures. Interests include multilingual/low-resource NLP, Cognitively-Inspired AI, Vision-Language Models and Adaptive AI for learners.

Experience

PicoLM

University of Cambridge *PicoLM, A Lightweight Framework for Studying Language Model Learning Dynamics*

📅 2024 Nov. – ongoing

📍 Cambridge, UK

- Supported by a generous research grant from the Accelerate Programme for Scientific Discovery at the University of Cambridge, I helped develop Pico, an open-source toolkit for pretraining small language models with systematic checkpointing of model states, activations, and gradients for tracking learning dynamics over time for developmental interpretability of Language Models Launched in March 2025, with over 50K+ views of launch video. Managing a team of Cambridge undergraduate students working on Pico in the next academic year.

Interactive Language Learning and Post-Training Researcher

Department of Computer Science & Technology *University of Cambridge*

📅 2025 March. – ongoing

📍 Cambridge, UK

- Leading a team of researchers from University of Cambridge, University of Sheffield and netmind.ai on Interactive Language Modelling and Emergent Communication.

BabyLM Workshop & Post-Training Research

In collaboration with KAIST AI, Naver & *BabyLM Workshop*

📅 2025 March. – ongoing

📍 Cambridge, UK

- Ongoing post-training work with KAIST AI group (XFACT, led by Dr James Thorne) and collaborators from University of Oxford, LinkedIn AI and Naver Cloud on Controlled Decoding and Process Reward Models.
- Pre-training research with Research Scientists from Naver Cloud, Korea and KAIST AI. Working on Multilingual BabyLM with Workshop Organisers.

Language Model Summer Researcher

Language Technology Lab & Department of Computer Science & Technology *University of Cambridge*

📅 2022 June. – 2022 September

📍 Cambridge, UK

- Extended research project with Prof Nigel Collier and Fangyu Liu (Google DeepMind) probing multimodal language models. Research Assistant and UROP Summer Intern roles on Argumentation Mining and Code-Switching.

Selected Publications

Pico. A Lightweight Framework for Studying Language Model Learning Dynamics

Diehl Martinez, R., Weiss, Yuval., Demitri Africa, David., SALHAN, Suchir., Daniels, Ryan & Buttery, Paula. ACL 2025 Systems Demonstration (Submitted)

Less is More. Pre-Training Cross-Lingual Small-Scale Language Models with Cognitively-Plausible Curriculum Learning Strategies.

SALHAN, S.A., Diehl Martinez, R., Goriely, Z., & Buttery, P. CoNLL, colocated in EMNLP 2024 (Miami, FL, USA)

Interests

Interactive Language Learning and Emergent Communication. Developmental Interpretability of Language Models.

Education

PhD in Computer Science, Cambridge University (2024 - 2028)

Guest Lecturer on Language Model Evaluation. MPhil Project Supervisor (Vision-Language Models). Teaching Assistant and Academic Supervisor Roles in NLP, Machine Learning & Bayesian Inference and Probability modules. Academic coursework in Geometric Deep Learning, Reinforcement Learning, Theory of Deep Learning & Mathematics of Machine Learning. Organiser of Departmental NLP Seminars.

BA, MEng in Computer Science & Linguistics, Cambridge University (2020 - 2024)

Starred First (Class I with Distinction) & Distinction. Modules include Advanced Topics in Machine Learning (XAI, Imitation Learning) and Deep Learning for NLP.

Skills

• Python • NumPy, SciPy, matplotlib, keras, scikit-learn, tensorflow • Java • React • Node.js • Docker • Git • HTML/CSS • MATLAB

Awards & Funding

Fully-Funded PhD Studentship

Memorial Prize

Memorial Prize, one of two highest academic accolades awarded to Gonville & Caius graduating cohort.

Senior Scholarship (Renewed) & Examination Prizes

Invited Talks

Human-Validated Grammar Profiles for Language Models

Tubingen (March 2025)