

# Engineering Small Language Models as Learner Models for LLM Interaction and Calibration

**Suchir Salhan**

Department of Computer Science & Technology, University of Cambridge, U.K. sas245@cam.ac.uk

## Small Language Models (SLMs) in ALTA Research: Goals and Rationale

Small Language Models (SLMs), typically under 1B parameters, provide valuable, interpretable, and efficient alternatives to large models. They are especially suitable for proprietary, task-specialised applications such as query routing in chatbot systems or edge/on-device ML.

## My PhD Work: SLMs as *Learner Models*

**Problem:** High-precision learner representations are essential for personalised and adaptive learning and assessment.

**Solution:** I propose bilingual SLMs for second language adaptation — or **L2LMs** — which simulate the developmental trajectories of second-language learners with a typologically-diverse L1s.

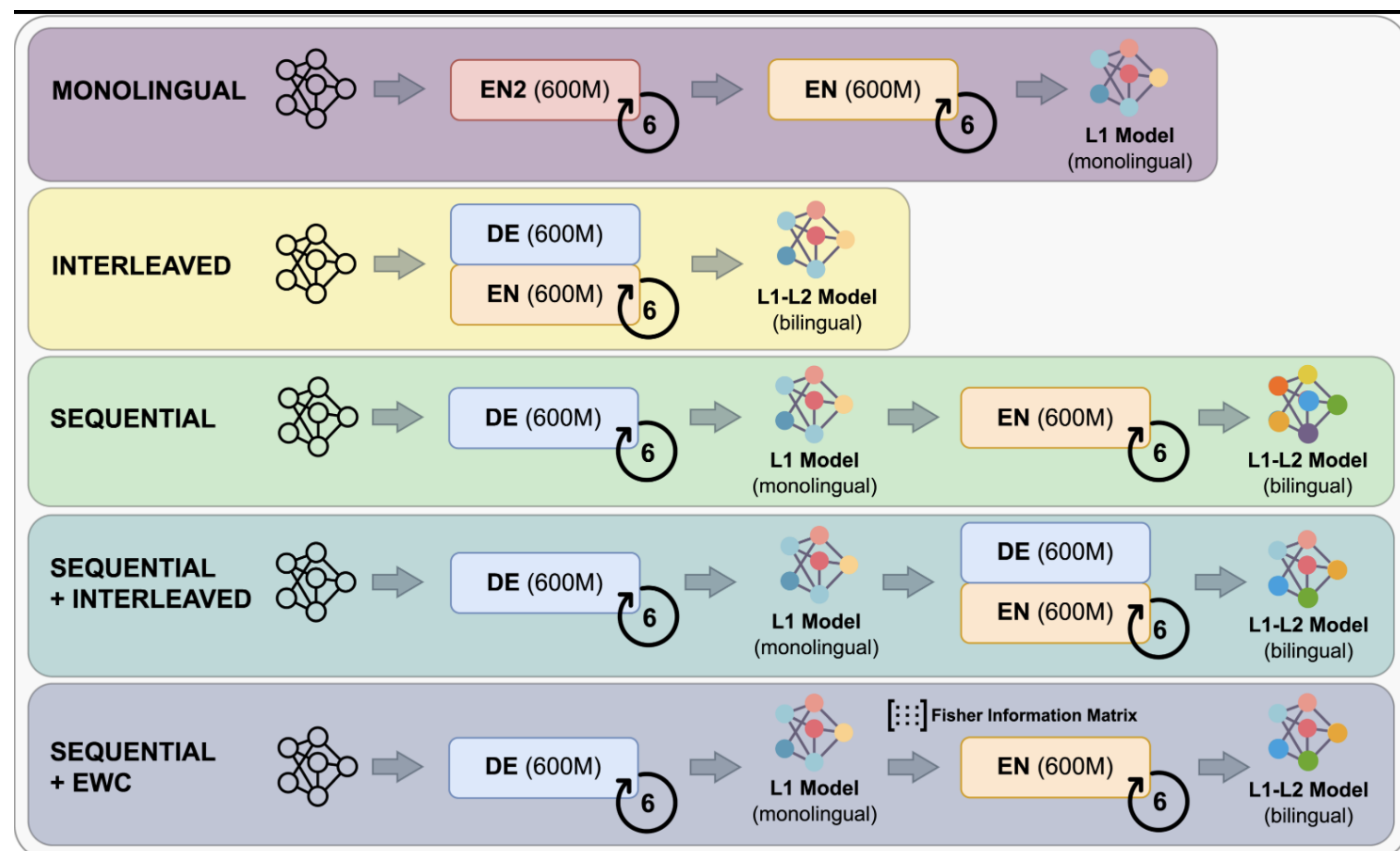


Figure 1: Conceptual Architecture of L2LMs and Sequential/Simultaneous Bilingual SLMs from Constantinescu et al (2025)

## SLMs and The North Star

**Design for EDIB:** Model design to reflect the *diversity of L1 backgrounds* trained on naturalistic volumes and distribution of learner corpora.

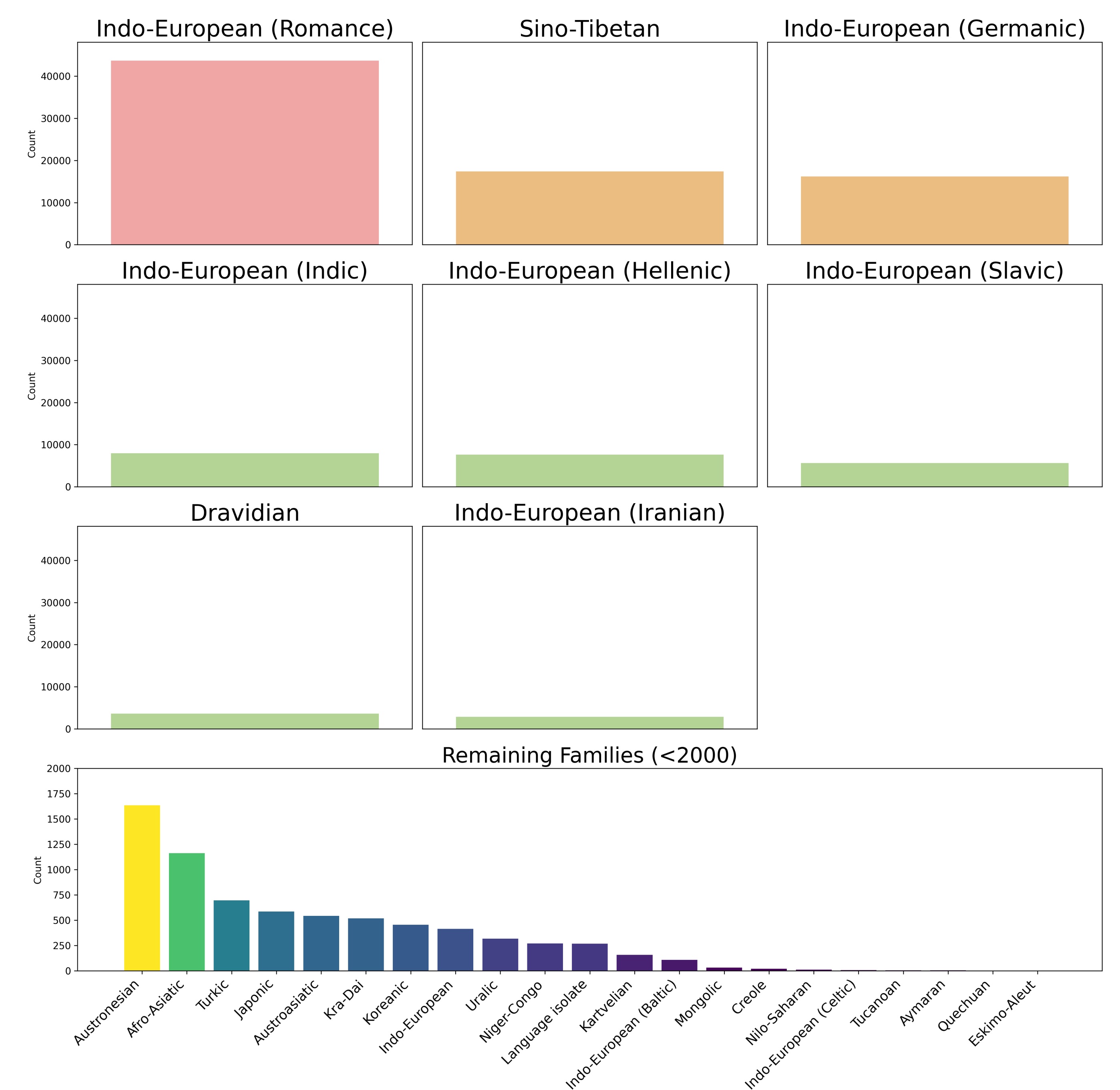


Figure 3: Distribution of L1 Backgrounds (by Language Family) in the Cambridge Learner Corpus

**#1 One Brand:** Integrating cutting-edge LLM techniques into L2LM pre-training [3,4] and post-training strategies for L2LM-LLM Interaction and Calibration.

**Trust through Interpretability:** Rich Learning Dynamics and Checkpointing of L2LMs in an explicit L2LM Design Philosophy, inspired by the Pico Learning Dynamics Framework [3].

## Selected References

- [1] Salhan, S.A (2025) **Linguistics in the Age of Language Models: What can Cognitively-Inspired Language Models offer to Linguistic Theory?** (Position Paper in *Cambridge Occasional Papers in Linguistics (CoPiL)*, Accepted, Volume 17) [2] Arnett, C. On the Acquisition of Shared Grammatical Representations in Bilingual Language Models *in press*
- [3] Diehl Martinez, R., Demitri Africa, D., Weiss, Y., Salhan, S.A., Daniels, R., & Buttery, P.J. (2025) **Pico: A Lightweight Framework for Studying Language Model Learning Dynamics** (under review)
- [4] Goriely, Z., Salhan, S.A., Lesci, P., Cheng, J., & Buttery, P.J. (2025) **ByteSpan: Information-Driven Subword Tokenisation** (Accepted ICML 2025 Tokenisation Workshop (TokShop, Non-Archival))
- [5] Salhan, S.A., Diehl Martinez, R., Goriely, Z., & Buttery, P.J. (2024) **Less is More: Pre-Training Cross-Lingual Small-Scale Language Models with Cognitively-Plausible Curriculum Learning Strategies** (Accepted Paper @ BabyLM Workshop in EMNLP 2024).

*This work is supported by Cambridge University Press & Assessment*

## Integration into ALTA CST

**Integrated Learning and Assessment:** Calibrating LLMs with L2LMs of L1 Background to *identify the right tests for the right purposes*. **Accessibility by Design:** Calibrations LLMs with L2LMs for lower-resourced L1s, especially those that occur infrequently in learner corpora, could help offer more *personalised learning profiles for a diverse population of English L2 learners*.