

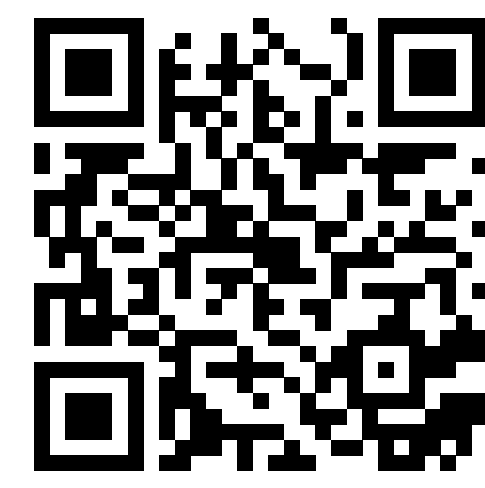
# Influence-driven Curriculum Learning for Pre-training on Limited Data

The First BabyLM Workshop @ EMNLP 2025

Loris Schoenegger<sup>1,2</sup>, Lukas Thoma<sup>1,2</sup>, Terra Blevins<sup>1,4</sup>, Benjamin Roth<sup>1,3</sup>

<sup>1</sup>Faculty of Computer Science, University of Vienna, Vienna, Austria

<sup>2</sup>UniVie Doctoral School Computer Science, University of Vienna, Vienna, Austria



universität  
wien

Faculty of Computer Science

<sup>3</sup>Faculty of Philological and Cultural Studies, University of Vienna, Vienna, Austria

<sup>4</sup>Khoury College of Computer Sciences, Northeastern University, Boston, USA

## Motivation

- Human-inspired CL performs poorly for low-resource pre-training
- Is there an inherent problem with sorting examples by difficulty?
  - Or are we just using the wrong heuristics?
  - Will a model-centric measure of difficulty perform better?

## Contributions

We leverage a technique from interpretability research to build a **novel type of curriculum**:

**Curricula based on training data influence estimates**

- We demonstrate their **effectiveness in benchmarks**;
- analyze their data mix** and how it evolves over time;
- study loss trajectories** to determine how they affect the model's learning process;
- and **compare example ordering** to existing sorting heuristics.

## Difficulty Score

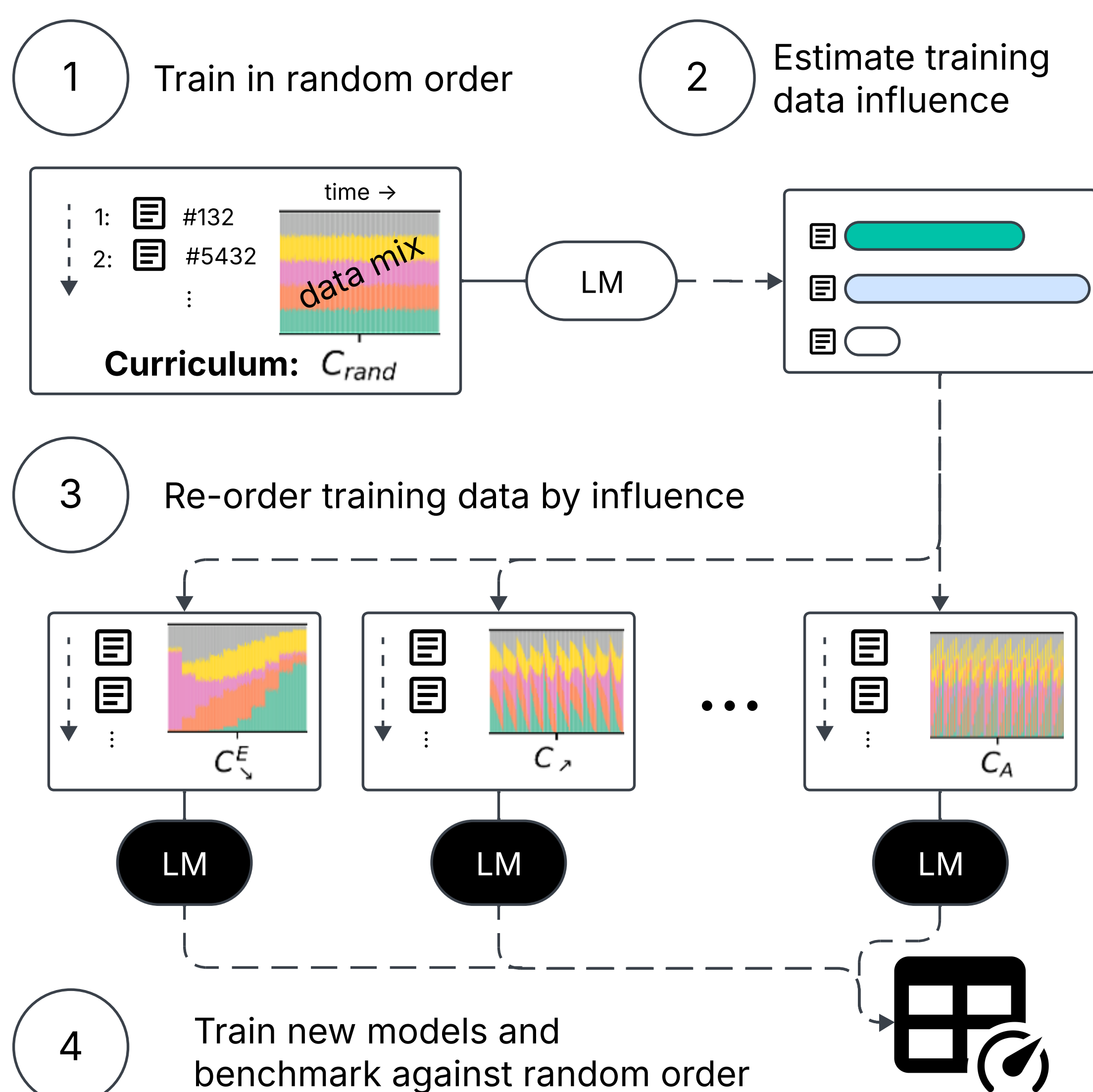
**Average influence**  $\phi_t(z, D)$  that a given training example exerts on the prediction of all other examples from the training data  $D$ :

$$\phi_t(z, D) = \frac{\sum_{z' \in D} \nabla \ell(w_t, z) \cdot \nabla \ell(w_t, z')}{|D|} \rightarrow (w_t \text{ is the WE layer})$$

$$= \nabla \ell(w_t, z) \cdot \mathbb{E}_{z' \sim D} [\nabla \ell(w_t, z')]$$

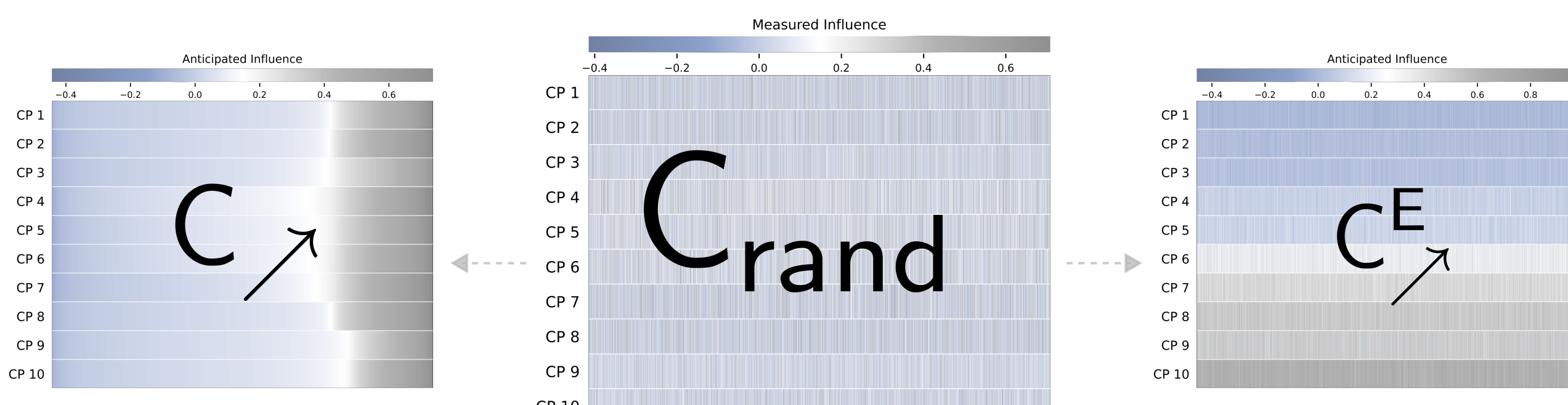
High for prototypical examples (whose loss gradients are similar to the average gradient), low for outliers.

## Curriculum Design



### Coverage Strategy Curriculum

Epoch-wise	$C_{\nearrow}$	Increasing difficulty, data visited once per epoch + 5 more
Cumulative	$C_{\nearrow}^E$	Easy examples in early-, difficult examples in later epochs + 2 more
Human-inspired	$C_{Source}$	Handcrafted source-difficulty curriculum (5 separate stages) + 2 more



## Datasets

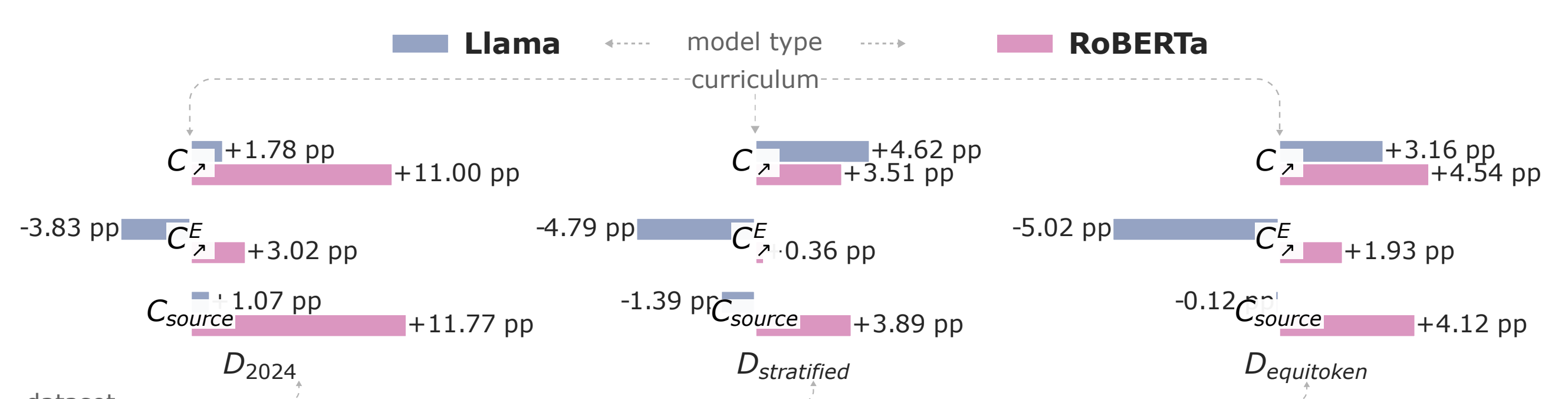
$D_{2024}$	2024/25 BabyLM dataset (10M word text-only)
$D_{stratified}$	equal number of words per stage
$D_{equitoken}$	equal number of words per example (100)

built from 5 stages:

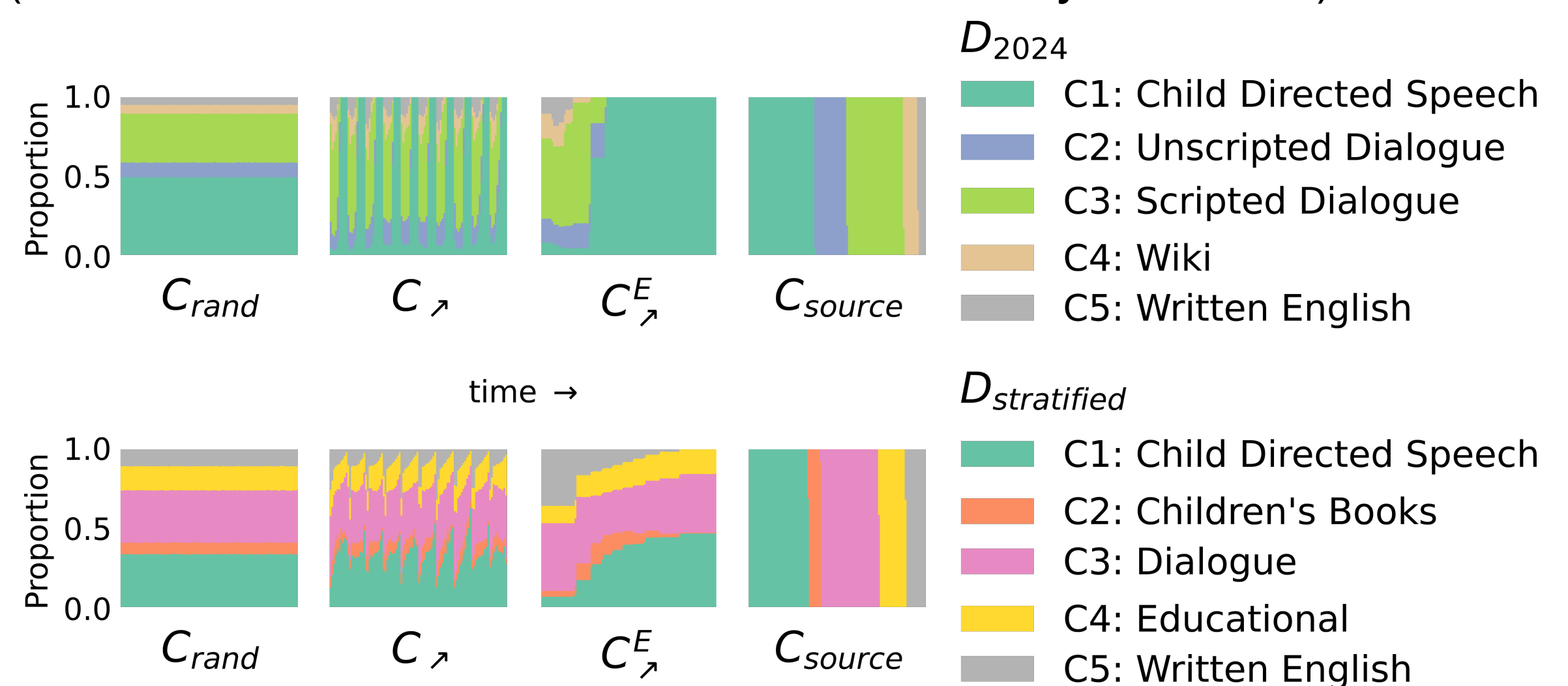
C1: Child Directed Speech  
C2: Unscripted Dialogue  
C3: Scripted Dialogue  
C4: Wiki  
C5: Written English

## Results

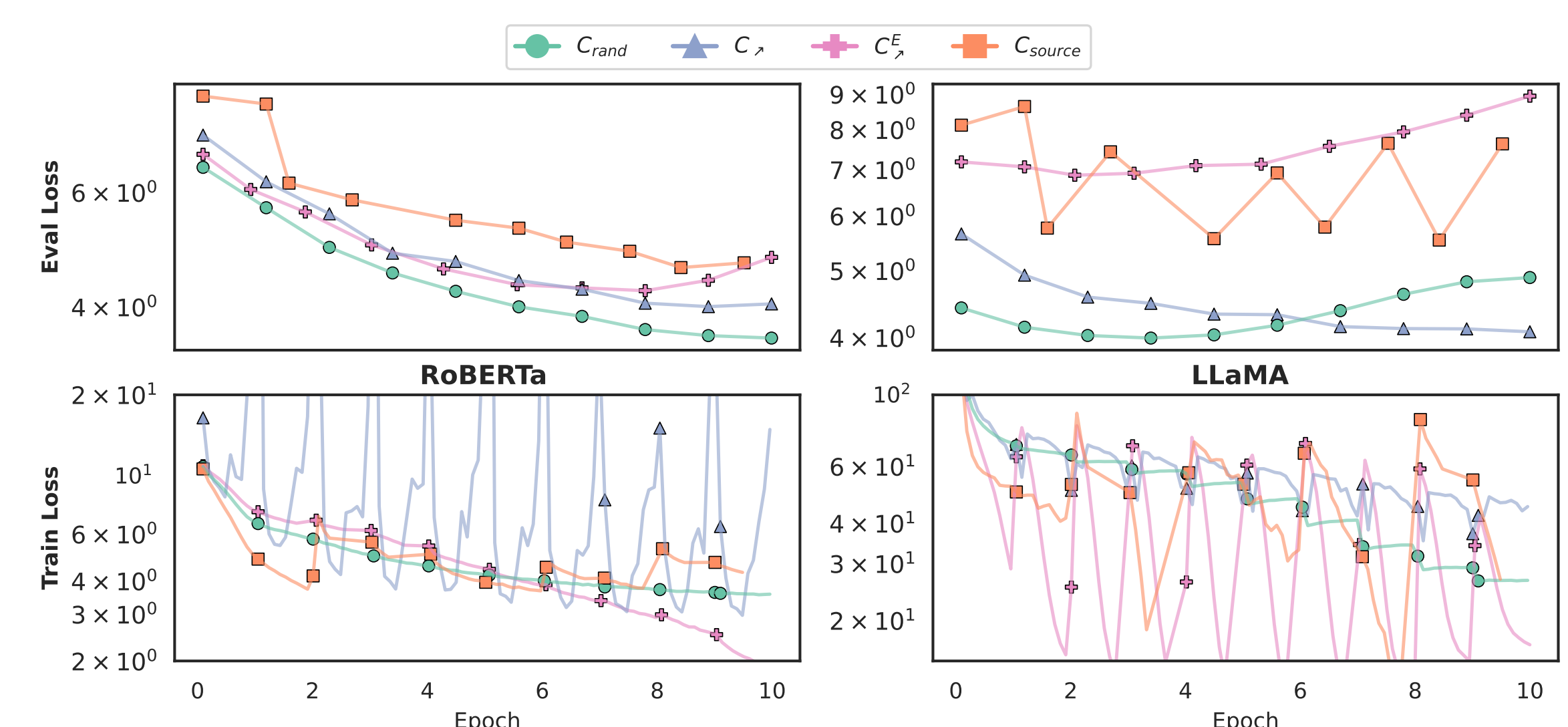
→ Increase of over 10 pp in acc over random order for RoBERTa- and over 4 pp for Llama models (see paper for best performing- and additional curricula):



→ source composition does not strongly vary over time (vs. traditional handcrafted source difficulty curricula):



→ severe spikes in training loss are not significantly correlated with performance on downstream benchmarks



## Shared Task Results

**Strict-Small Track.**  $E$  denotes BabyLM baseline models.

Model	GLUE	blimp.filt	blimp.supplement	comps	entity_tracking	ewok	Avg acc
$TICL_{roberta}$	0.646	0.699	0.539	0.507	0.343	<b>0.503</b>	0.598
$E_{masked}$	<b>0.665</b>	0.502	0.480	0.500	<b>0.422</b>	<b>0.503</b>	0.496
$E_{mixed}$	0.660	0.501	0.467	0.491	0.414	0.500	0.493
$E_{causal}$	0.654	<b>0.717</b>	<b>0.632</b>	<b>0.528</b>	0.346	0.495	<b>0.614</b>
$E_{opt2}$	0.623	0.664	0.571	0.517	0.139	0.499	0.540

Correlation in wug tasks

(Spearman's  $\rho$  Hofmann et al., 2025).

Model	wug.adj.nom	wug.past.tense	Avg
$TICL_{roberta}$	0.006	-0.001	0.003
$E_{masked}$	0.005	-0.002	0.001
$E_{mixed}$	0.004	-0.001	0.002
$E_{causal}$	0.006	<b>0.001</b>	<b>0.004</b>
$E_{opt2}$	<b>0.007</b>	-0.001	0.003

Reading tasks. Reported as %  $R^2$  gain.

Model	Eye Tracking Score	Self-Paced Reading Score	Avg
$TICL_{roberta}$	0.040	0.002	0.021
$E_{masked}$	<b>0.103</b>	0.027	0.065
$E_{mixed}$	0.099	0.025	0.062
$E_{causal}$	0.099	0.035	<b>0.067</b>
$E_{opt2}$	0.087	<b>0.043</b>	0.065

## Key Findings

- Our sorting strategies can increase performance; however: **only if paired with non-developmentally plausible dataset coverage strategies**, i.e., must **visit the full dataset every epoch**
- Improvement may result from improved grouping of examples into **batches of similar difficulty**
- Measure appears **inversely correlated** to those of other sorting heuristics (high influence → low difficulty)