

BabyBabelLM: A Multilingual Benchmark of Developmentally Plausible Training Data

Anonymous ACL submission

Abstract

We present **BabyBabelLM**, a multilingual collection of datasets modeling the language a person observes from birth until they acquire a native language. We curate developmentally plausible pretraining data in each of 45 languages, aiming to cover the equivalent of 100M English words of content. We compile evaluation suites and train baseline models in each language. BabyBabelLM aims to facilitate multilingual pretraining and cognitive modeling.¹

1 Introduction

The prevailing trend in language modeling research is to prioritize scaling, both in terms of model size and training data volume (Kaplan et al., 2020; Choshen et al., 2024). While this approach has led to significant advances in model performance, it neglects fundamental research questions about the nature of language learning (Wilcox et al., 2024). It disincentivizes work on data-efficient modeling, which, from a practical perspective, offers benefits in terms of efficiency and accessibility. From a theoretical perspective, it ignores the growing mismatch between human language acquisition and language model (LM) learning. From infancy to maturity, English learners acquire language through exposure to less than 100M words (Gilkerson et al., 2017), several orders of magnitude less than the massive pretraining corpora required by contemporary LMs surpassing 10T words (Bengio et al., 2025).

In response to the field’s focus on scale, the BabyLM Challenge (Warstadt et al., 2023) was created to redirect attention toward questions of data efficiency and developmental plausibility in language modeling. The shared task invites participants to propose data-efficient LMs pretrained on a fixed, developmentally plausible English corpus

of child-directed speech (CDS), educational content, and other simplified texts. The top-performing submissions (Charpentier and Samuel, 2023, 2024) have significantly improved the state of the art for models trained on the same limited data budget, even surpassing LMs trained on much larger corpora on various benchmarks.

The BabyLM Challenge has generated a new line of research on data-efficient training and cognitively-inspired modeling (Warstadt et al., 2023; Hu et al., 2024), depending on the existence of developmentally plausible datasets as training corpora and supplying resources to ease and direct such research. However, the majority of this work has focused on English, largely due to the public availability of the pretraining corpora released for the BabyLM challenge, which is English-only. There is a small but growing body of work that extends the BabyLM research project beyond English (Salhan et al., 2024; Shen et al., 2024; Prévot et al., 2024; Matzopoulos et al., 2025; Padovani et al., 2025; Bunzeck et al., 2025). Such efforts are crucial for developing an accurate understanding of the relationship between human language acquisition and LM learning. Any claim that a technique is developmentally plausible can only be truly substantiated by evaluations across typologically diverse languages. In addition, there is variation in language acquisition between languages, and human language learning also occurs in multilingual settings (Grosjean, 1989; Slobin, 2014).

To facilitate this research, we create BabyBabelLM, a multilingual collection of developmentally plausible training datasets. The collection includes 45 languages, encompassing families primarily rooted—though not exclusively spoken—in Europe, Asia, and Africa. For each language, we carefully select and compile publicly available datasets, as well as release new ones. This includes several categories of developmentally plausible data, such as CDS, educational resources, and

¹All code and data will be released upon publication.

other child-oriented content (e.g., books, news, and wikis aimed towards children). As the languages in BabyBabelLM vary widely in terms of public availability of developmentally plausible data, the exact composition of the datasets is not identical, though we prioritize plausible data sources in all cases. We sort languages into three tiers based on training set size, corresponding to the equivalent of respectively 100M, 10M, or 1M English words, calibrated by language-adjusted byte estimates (Arnett et al., 2024), to ensure comparability of data budgets across languages with differing orthographic and morphological characteristics.

To further facilitate research we also compile a list of evaluations to test models created on those domains. While we focus on linguistic aspects that are more clearly related to language learning rather than domains of the training data, we provide a comprehensive list of existing datasets to facilitate any future questions and to provide coverage. Specifically, the datasets allow evaluation that fits the pretraining objective directly, with no dataset-specific adaptation (known as zero-shot).

Overall, this effort releases:

- Developmentally plausible **pretraining datasets** for 45 languages, released under licenses permitting research purposes (§3).
- A **pipeline** to allow for subsequent dataset expansion with new resources and languages (§3.3).
- A survey of multilingual **evaluation** tasks (§4) accompanied by an evaluation suite extendable by the community.
- A collection of 45 monolingual **pretrained models**, 7 bilingual models and a multilingual model that we analyze in §5.

2 Related Work

The first edition of the BabyLM challenge (Warstadt et al., 2023) released two pretraining corpora, respectively 10M and 100M words, each consisting of 39% developmentally plausible data and a selection of high-quality corpora (e.g. Wikipedia). The second edition (Hu et al., 2024) updated the datasets to increase the proportion of child-oriented data to 70%. Thus far, the BabyLM Challenge has been limited to English for training and evaluation. In both editions, BabyLM submissions were evaluated on two types of language tasks: Minimal pair challenges (Warstadt et al., 2020; Ivanova et al., 2024) testing linguistic competence, world knowledge or other capabilities by asserting the

model prefers a correct sentence over a sentence with a minor but meaningful alteration. In contrast, other types of evaluations (which we exclude in this work) fine-tune the model on a dataset and test its ability to learn a new task.

Beyond English, a growing body of work has begun exploring BabyLM-style models and the collection of developmentally plausible training datasets for other languages. By extending BabyLM research to typologically diverse languages, these studies have produced new insights into cognitively-inspired language modeling. Salhan et al. (2024) propose acquisition-inspired curriculum learning strategies and train small-scale LMs on age-ordered CDS for French, German, Japanese, and Chinese. Prévot et al. (2024) investigate the value of spontaneous speech corpora for BabyLM evaluation with experiments on English and French. Matzopoulos et al. (2025) train BabyLMs for isiXhosa, a low-resource South African language, highlighting the limits of BabyLM research for languages without publicly available developmentally plausible corpora. Capone et al. (2024) release a corpus of Italian developmentally plausible training data. Padovani et al. (2025) show that training on CDS does not consistently improve grammatical learning across English, French, and German. Bunzeck et al. (2025) train LMs on distributionally varied subsets of a German BabyLM corpus, showing that syntax learning benefits from complex constructions while lexical learning benefits from fragmentary constructions. Finally, Shen et al. (2024) investigate developmentally plausible L2 acquisition by adapting an English BabyLM for Italian via a reward signal from a parent Italian model. However, these works are typically forced to compile novel datasets in addition to their scientific contribution, and they do not represent a coordinated effort to compile such training data in comparable ways and across diverse languages.

Some relevant multilingual resources do exist. The Child Language Data Exchange System (CHILDES; MacWhinney, 2000) is a multilingual database of transcriptions of child-adult interactions, including data for over 40 languages, with varying child age ranges, interaction environments, and corpus sizes. CHILDES serves as a starting point for most of our languages in BabyBabelLM. A previous effort to compile developmentally plausible multilingual training corpora is MAO-CHILDES (Yadavalli et al., 2023), an age-

ordered dataset of CHILDES corpora for five typologically diverse languages (German, French, Polish, Indonesian, and Japanese), which the authors use to study cross-lingual training and L2 learning. [Salhan et al. \(2024\)](#) and [Goriely and Buttery \(2025\)](#) independently release MAO-CHILDES and IPA-CHILDES for four languages (Japanese, Chinese, French, German) and a phonemized corpus based on CHILDES for 31 languages.

3 Dataset Creation and Overview

The BabyBabelLM dataset was constructed to support research on developmentally plausible language modeling across a wide range of languages. Our aim is to approximate the kind of linguistic input that humans are exposed to in early life, while providing clean, well-documented, high-quality data. This section outlines how the dataset was created: the principles of data selection, the structure and content of the data, and the data preprocessing pipeline.

3.1 Data Collection Principles

The design of our datasets required various methodological choices regarding the types of data sources to include, ensuring their developmental plausibility, and ensuring long-term extensibility. In this section, we describe the criteria guiding our choices, the organizational structure of our multilingual collection, and the licensing considerations.

3.1.1 Developmental Plausibility Criteria

Our guiding principle in dataset construction is that of **developmental plausibility**: the idea that pretraining data should approximate the linguistic input children encounter. To this end, we prioritized domains such as child-directed speech (CDS), educational materials, children’s books, and transcribed conversations. We also made a deliberate choice to exclude synthetic corpora, such as TinyStories ([Eldan and Li, 2023](#)) or TinyDialogues ([Feng et al., 2024](#)), despite their developmental intention. Manual inspection of these datasets revealed that they often contain unrealistically verbose or structurally complex language. In addition to content filtering, we prioritized data quality by removing noise, favoring conversational data when applicable, and standardizing the format of metadata (see Appendix A). This preserved realism enables controlled cross-lingual comparisons, which are essential for studying the impact of linguistic variation on model learning.

3.1.2 Community-driven Data Leadership

To ensure dataset quality, data collection for most languages was led by a researcher fluent in or highly familiar with that language. These language leads were responsible for sourcing appropriate corpora, verifying developmental plausibility, and coordinating with local experts in linguistics and language acquisition.

The BabyBabelLM dataset is designed as a “living resource”. As more developmentally plausible data becomes available, we aim to expand the collection both in breadth—by adding new languages—and in depth—by enriching existing ones. To support this, we provide an open-source pipeline that enables researchers to add entirely new languages and expand existing language datasets. While our initial release covers 45 languages, 16 languages rely solely on general-purpose multilingual data resources. We consider these entries as starting points for future, more comprehensive corpora.

We invite contributions through GitHub and Hugging Face, where researchers can submit new datasets, improvements, and evaluations. All additions are reviewed for compliance with our guidelines and incorporated into future versions of the dataset, ensuring proper attribution. We hope this model of open, collaborative development will lead to broader coverage and increased utility across diverse research agendas.

3.1.3 Licensing and Ethics

During our data collection effort, we verified that all data is released with licenses that permit academic research, such as Creative Commons or Public Domain. When licensing information was missing, the right holders of each source were contacted. We release our corpus with document-level licensing information and data source attribution to ensure that each resource is used ethically and within its rights. In the rare cases where no license or contact information was available, we decided to still release the data, but under a restrictive non-commercial license.

3.2 Dataset Composition

Constructing a multilingual dataset that is both developmentally plausible and broadly comparable requires careful attention to how data is organized. Languages differ widely in the availability and type of child-relevant resources, and these differences must be accounted for without undermining cross-

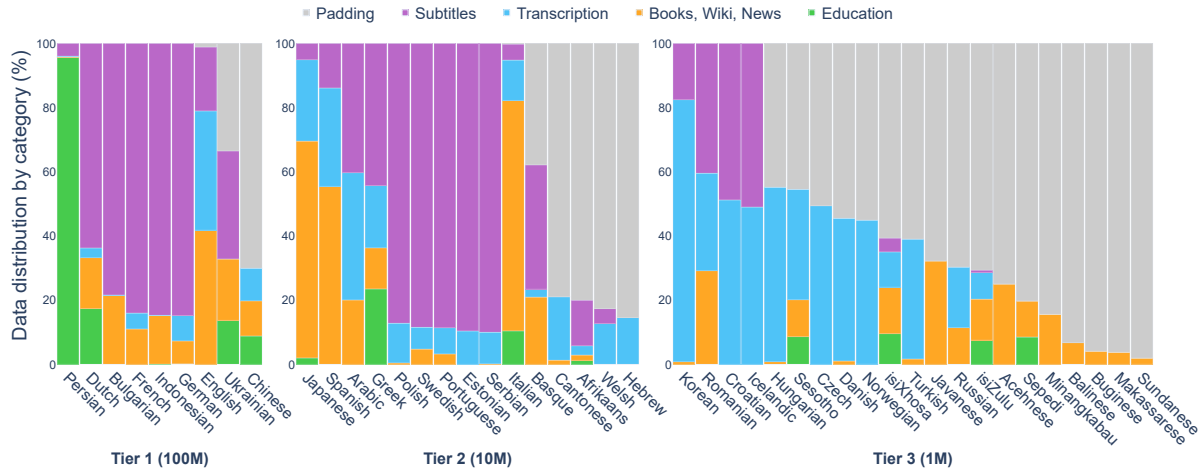


Figure 1: Training data distribution by category across languages for all data tiers in the BabyBabelLM dataset.

linguistic consistency. This section outlines how we approached this challenge, the types of data we included, and how they differ from one another.

3.2.1 Data categories

Transcription Children learn language primarily from spoken input, which we therefore use as our primary data source. This child-directed speech (or *CDS*) differs drastically from the language data found in commonly used pretraining corpora. Usually, it is structurally short and simple (Genovese et al., 2020), and features high amounts of syntactic and lexical repetition (Tal et al., 2024), while its vocabulary is mostly restricted to everyday topics and children’s immediate surroundings (Snow and Ferguson, 1977). The CHILDES database (MacWhinney, 2000) contains a large amount of such data in the form of recorded caretaker-child interactions (e.g., during free play, meal times, or shared book reading) and manually created transcriptions. We used all CDS available for our target languages in CHILDES as the base of our datasets.² As children also overhear language in their surroundings, we further included as much child-available speech (adult-adult dialogue) as possible per language.

Education We included educational content aimed at children, taken from textbooks and exams. As children spend a large amount of their time in education systems, they encounter this kind of input regularly. On the content level, it provides

much more direct instruction than CDS, which we deem useful for our purposes. After all, our BabyBabelLMs are not only supposed to learn formal linguistic patterns from the input, but ideally also more functional (visual semantic, pragmatic, and world) knowledge.

Books, Wiki, News To approximate the whole breadth of input that children receive, we further included child-oriented media, i.e., children’s books, children’s wikis, child-targeted news, and other appropriate media sources. For multilingual resources, we incorporated the Ririro story collection³, GlotStoryBooks from Kargaran et al. (2023), and Child Wiki articles across many languages. Additionally, individual languages were enriched with monolingual resources. In contrast to CDS, this kind of data features longer and more complex sentences (Cameron-Faulkner and Noble, 2013; Bunzeck et al., 2025), and much more diverse vocabulary and content. As such, these sources should provide a useful training signal for more complex knowledge levels, similar to educational content.

Subtitles Finally, we also used movie/TV show subtitles from OpenSubtitles (Lison and Tiedemann, 2016). While such fictional speech does differ from natural spoken data – for example, it features less hesitations, interjections, false starts or pauses (Bishop, 1991; Jucker, 2021; Gast et al., 2023) – it still approximates the linguistic properties of speech well and we deem it developmentally plausible. Furthermore, children are nowadays exposed to a wide variety of (video) media (Gowen-

²For some languages not written in Latin script (e.g., Japanese, Greek, Persian), the CHILDES data contains transliterations. We include this data *as-is* to maintain comparability with existing resources, noting that back-transliteration introduces ambiguity and may reduce data quality.

³<https://ririro.com/>

lock et al., 2024), and thus encounter this kind of content regularly. We did omit certain categories (e.g., adult content, crime, horror) to ensure our datasets do not contain content inappropriate for children. As the OpenSubtitles corpora are significantly larger than our other data sources, we did not include all of their content; instead, we used them to pad our datasets to the different tiers.

3.2.2 Language Tiers and Coverage

Our dataset spans 45 languages drawn from a wide range of typological families. While Indo-European languages are well represented (22 out of 45), the collection also includes Semitic, Uralic, Bantu, Austronesian, and Sino-Tibetan languages, among others. This diversity was a key design goal, enabling investigation of language acquisition and model performance across structurally distinct linguistic systems.

However, linguistic diversity is closely tied to disparities in data availability, resulting in big variations in data quantities for our set of languages. To enable fair comparisons, we classify languages into three distinct **tiers** according to the amount of collected data. Tier 1 includes languages with roughly 100 million English-equivalent tokens, Tier 2 with 10 million, and Tier 3 with 1 million. Ranking them by decreasing dataset size, the tiers contain 9, 15, and 21 languages respectively. This distribution further underscores the current scarcity of developmentally plausible corpora and the need for community-driven collection efforts.

Token thresholds are *calibrated* using the **byte premium** approach (Arnett et al., 2024), which adjusts for variation in orthographic and morphological structure by measuring the UTF-8 encoded size needed to express a fixed amount of content. For each language, we curated as much developmentally plausible content as possible before padding to the tier threshold using fallback data sources such as OpenSubtitles. When subtitle data was insufficient, additional fallback corpora—including FineWeb-C and Wikipedia—were used, with filtering to ensure appropriateness. Figure 1 summarizes the distribution of content categories across languages and tiers; more detailed per-language statistics are presented in Table 2.

While our data selection prioritizes developmental plausibility, public availability varies widely, and not all languages include the same balance of content types. We caution that downstream comparisons across languages may be influenced by

these differences in dataset composition.

3.3 Data Preprocessing

The data preprocessing is separated into two stages. Initially, language-specific preprocessing was carried out by the language leads, as needed by the specific data and language (more in Appendix B).

Following language-specific preprocessing, we apply (and release) a uniform pipeline to all data, including standard normalization (unicode, whitespace, punctuation) and category-specific preprocessing. For dialogue transcripts, we remove linguistic annotations while preserving tab-separated utterance structures. For subtitle data, we remove speaker labels, music note symbols, stage directions, and timestamps. For book-like formats (educational materials, children’s books, wikis) and the QED dataset, we remove XML tags and URLs.

For language and script validation, we use GlotLID v3 (Kargaran et al., 2023) to classify sentence-like chunks of text. These chunks are created by first splitting documents into paragraphs followed by sentence-based heuristics. The document’s final language is assigned via a segment-based majority vote. To maintain data quality, we filter mismatched segments within documents and discard any document that fails the overall validation. Other document metadata fields, such as the text category and license, are validated for the correct type and values when applicable (see Table 3).

4 Evaluation Suite

We create a multilingual evaluation suite that targets both the *formal linguistic competence* (knowledge of linguistic rules and patterns), and the *functional linguistic competence* (understanding and using language in the world) (Mahowald et al., 2024). We reviewed a large number of existing multilingual and monolingual benchmarks (Huang et al., 2025) with the aim of ensuring all our languages have at least one evaluation dataset testing formal and one testing functional linguistic competence.

Formal competence To assess formal linguistic competence, we prioritized high-quality, language-specific minimal pair benchmarks that target a diverse set of linguistic phenomena. This approach was applied to languages such as Basque (Kryvosheieva and Levy, 2025), Chinese (Liu et al., 2024), Japanese (Someya and Os-eki, 2023), German (Vamvas and Sennrich, 2021),

and Turkish (Başar et al., 2025). For other languages, we employed datasets that cover fewer phenomena but span a broader range of languages. In particular, for English, French, German, Russian, and Hebrew, we used CLAMS (Mueller et al., 2020), a cross-lingual minimal pair benchmark built from linguist-curated templates, focusing on subject-verb number agreement across various syntactic contexts. In our experiments we refer to the collection of these tasks as *MonoBLiMP*. Finally, we incorporated MultiBLiMP (Jumelet et al., 2025), a large-scale dataset of minimal pairs automatically generated from the Universal Dependencies treebanks (Nivre et al., 2017). MultiBLiMP targets subject-verb agreement in number, person, and gender, and offers the widest language coverage among our benchmarks, as detailed in Table 1.

Functional competence We include two types of benchmarks to evaluate functional competence. The first category focuses on factual and domain-specific knowledge memorized by the model, such as Global-MMLU (Singh et al., 2025), INCLUDE (Romanou et al., 2024), and BM-LAMA (Qi et al., 2023). The second category assesses general reasoning abilities, including natural language inference, commonsense reasoning, narrative understanding, and reading comprehension. Benchmarks in this category include XNLI (Conneau et al., 2018), MultiNLI (Williams et al., 2018), HellaSwag (Zellers et al., 2019), Belebele (Bardkar et al., 2024), ARC (Clark et al., 2018), xstorycloze (Lin et al., 2022b), TruthfulQA (Lin et al., 2022a), XCOPA (Ponti et al., 2020), SIB-200 (Adelani et al., 2024), and XWinograde (ai2, 2019; Cheng and Amiri, 2024). Additionally, we included XCOMPS (He et al., 2025), a multilingual conceptual minimal pair dataset with 17 languages.

Evaluation For conducting these evaluations, we relied on Eleuther AI’s LM Evaluation Harness (Gao et al., 2024). All our evaluations were framed as zero-shot multiple-choice problems where models’ answers are determined by selecting the option with the highest log probability. We used reference implementations for some tasks and implemented the ones that were missing.

Eight of our languages remain very low-resourced in terms of evaluation (i.e., evaluated only by one NLU task). These languages are Welsh, Yue Chinese, Achinese, Balinese, Buginese, Croatian, Makasar, and Minangkabau. Developing BabyLM-suitable benchmarks for these and other

languages remains an important priority for future work. To facilitate research in this direction, we publicly release our evaluation suite and plan to support community contributions (see §3.1.2).

5 Experiments

Building on the resources outlined above, we train monolingual, bilingual and multilingual models to evaluate our benchmark suite. These models serve as simple, reproducible baselines, providing a foundation for a broader, community-driven challenge.

Setup For training our models, we adopt the model configurations of the GoldFish model suite (Chang et al., 2024). For the monolingual models, we use a lightweight GPT-2 architecture with 4 transformer layers, 8 attention heads, and a hidden size of 512. The model uses GELU activations and standard dropout regularization (0.1) across attention, embeddings, and residual connections. It includes a feedforward inner dimension of 2048 and supports sequences up to 512 tokens. For all languages, we use a BPE tokenization (trained on the training corpus), with a vocabulary size of 8,192 tokens (Huebner et al., 2021). This results in small LMs of only 17.1M parameters. Each model is trained for 10 epochs.

For the bilingual models, we train a model using data from each language in Tier 1 and the English BabyLM (200M tokens total), keeping model configuration the same, but reducing training epochs to 5. For the multilingual model, we increase the number of layers to 12, hidden size to 768, and vocabulary size to 32,768, accommodating the wide range of languages and scripts this model should handle. The model is trained for only 1 epoch, and has 111M parameters. We additionally compare performance against Qwen3-0.6B (Yang et al., 2025), a capable multilingual LM of modest scale.

Results The results for our monolingual models are presented in Table 1. Linguistic benchmarks such as MultiBLiMP yield promising results, with Tier 1 models typically scoring above 80%. Performance on MultiBLiMP is strongly driven by data size, with Tier 2 and 3 languages performing worse. Performance on other benchmarks remains close to random chance (e.g., XCOPA, ARC, XCOMPS, HellaSwag). As such, our comparatively tiny BabyLMs only provide a starting point for further experimentation.

We further compare the results of our mono-

Language	multiblmp	monoblmp	bmlama	xcopa	arc	xcomps	hellaswag	xnli	multinli	wino grande	sib200	belebele	global-mmlu	include	xstorycloze	truthfulqa
Random	50.0	50.0	10.0	50.0	25.0	50.0	25.0	33.3	33.3	50.0	14.3	25.0	25.0	25.0	50.0	25.0
TIER 1 (100M)																
Bulgarian	90.8	—	24.7	—	25.3	—	26.8	34.5	32.0	52.0	19.6	28.3	—	28.0	49.5	25.5
Chinese	—	70.2	34.5	51.2	36.9	55.1	26.8	33.3	35.8	49.2	11.3	22.3	23.0	23.9	48.7	22.1
Dutch	90.5	—	29.4	—	30.3	52.4	26.5	—	34.6	50.0	9.3	27.0	24.5	32.4	49.1	22.4
English	82.0	65.9	31.8	58.0	28.8	—	26.5	36.3	31.8	51.4	24.0	23.0	23.2	—	49.5	22.4
French	94.1	69.7	22.7	—	27.5	50.6	26.4	36.1	35.7	51.3	8.8	27.3	23.9	22.7	47.6	20.4
German	88.6	77.1	26.3	—	27.1	52.6	25.9	34.9	32.3	51.8	11.8	24.7	24.6	16.3	48.8	23.8
Indonesian	—	—	25.6	53.6	30.4	—	27.3	—	31.7	53.1	10.8	25.4	26.1	22.9	50.5	24.3
Persian	71.3	—	29.6	—	29.7	53.6	26.4	—	35.3	50.7	9.8	26.0	23.8	23.0	52.2	26.9
Ukrainian	88.6	—	24.8	—	28.4	50.6	26.4	—	31.6	50.3	8.3	22.8	23.4	30.6	47.6	24.7
TIER 2 (10M)																
Afrikaans	—	—	28.5	—	28.3	—	26.1	—	32.0	51.3	10.8	22.1	—	—	49.5	22.6
Arabic	75.9	—	17.2	—	25.7	52.9	25.8	33.0	32.4	47.4	19.1	22.9	23.1	22.1	46.8	22.3
Basque	94.5	65.3	—	49.8	28.5	—	—	33.6	—	—	10.3	26.0	—	25.4	50.6	28.0
Estonian	81.5	—	21.0	53.0	24.8	—	25.5	—	35.9	50.7	12.2	22.0	—	21.4	45.8	22.6
Greek	89.2	—	23.2	—	26.9	50.3	26.4	35.6	31.7	49.5	19.6	22.8	23.1	20.9	49.4	25.7
Hebrew	70.2	59.5	23.0	—	29.8	51.6	26.1	—	32.5	50.2	12.2	22.9	23.1	24.8	49.6	23.1
Italian	77.5	—	26.3	52.0	26.4	—	26.5	—	31.7	50.1	11.3	24.6	23.3	20.7	50.2	25.2
Japanese	—	61.9	19.9	—	27.8	50.8	25.1	—	31.7	47.5	19.6	23.4	23.0	26.1	47.8	22.1
Polish	75.9	—	19.4	—	26.1	—	25.5	—	32.0	49.0	12.2	21.9	26.6	18.2	49.5	21.8
Portuguese	80.7	—	21.7	—	25.4	—	26.3	—	31.7	48.8	12.2	23.9	23.5	22.2	48.8	23.6
Serbian	—	—	24.4	—	25.5	—	25.6	—	32.4	49.5	8.3	22.9	23.1	20.2	48.5	23.5
Spanish	83.0	—	24.1	—	28.0	51.3	26.4	35.5	33.7	49.1	18.6	24.9	24.4	23.7	47.8	23.0
Swedish	—	—	23.3	—	26.2	—	25.9	—	31.7	49.3	12.2	28.3	24.6	—	48.1	24.5
Welsh	91.4	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
Yue Chinese	—	—	—	—	—	—	—	—	—	—	19.6	—	—	—	—	—
TIER 3 (1M)																
Achinese	—	—	—	—	—	—	—	—	—	—	18.1	—	—	—	—	—
Balinese	—	—	—	—	—	—	—	—	—	—	23.5	—	—	—	—	—
Buginese	—	—	—	—	—	—	—	—	—	—	19.6	—	—	—	—	—
Croatian	—	—	17.8	—	26.1	—	25.8	—	31.7	47.6	12.2	22.4	—	25.4	45.6	21.1
Czech	59.0	—	13.1	—	24.8	—	25.8	—	35.9	50.4	19.6	22.9	23.1	—	48.6	20.2
Danish	82.0	—	17.2	—	24.6	—	25.7	—	35.9	49.0	12.2	23.6	—	—	48.9	21.8
Hungarian	68.9	—	10.2	—	24.3	49.9	25.6	—	31.7	48.0	12.2	22.9	—	30.4	48.5	18.9
Icelandic	71.6	—	13.4	—	24.3	—	25.4	—	31.7	50.0	12.2	22.9	—	—	45.7	22.6
Javanese	—	—	17.4	—	25.2	—	25.8	—	31.7	50.5	19.6	28.2	—	—	50.6	24.1
Korean	—	—	19.1	—	25.5	51.3	25.0	—	32.5	48.9	10.8	21.6	26.5	21.8	47.5	23.3
Makasar	—	—	—	—	—	—	—	—	—	—	12.4	—	—	—	—	—
Minangkabau	—	—	—	—	—	—	—	—	—	—	19.6	—	—	—	—	—
Norwegian	—	—	21.4	—	26.4	—	25.9	—	35.9	47.0	12.2	22.9	—	—	47.1	22.3
Pedi	—	—	—	—	25.9	—	—	—	—	—	10.8	24.1	—	—	—	—
Romanian	74.1	—	15.4	—	25.3	—	25.5	—	35.7	48.6	14.7	24.0	24.4	—	46.2	21.8
Russian	58.5	52.0	17.2	—	25.2	49.1	25.9	33.3	35.9	51.3	12.2	28.9	24.9	20.5	48.8	21.4
Southern Sotho	—	—	—	—	—	—	—	33.3	—	—	9.3	22.9	20.0	—	—	—
Sundanese	—	—	—	—	—	—	—	—	—	—	19.6	27.2	—	—	—	—
Turkish	64.9	59.8	17.9	53.8	25.4	51.6	25.9	33.3	35.9	49.9	12.2	21.9	27.1	23.0	50.1	24.7
Xhosa	—	—	—	—	—	—	—	33.3	—	—	10.8	22.9	20.0	—	—	—
Zulu	—	—	—	—	29.7	—	—	33.3	—	—	10.8	23.0	22.8	—	—	—

Table 1: Performance of the monolingual models trained on BabyBabelLM. All scores denote average 0-shot accuracy. Columns are sorted by difference of the average task performance against random chance.

lingual models for MultiBLiMP and Belebele with the multilingual BabyLM model (Multi-BabyBabelLM) and the Qwen3-0.6B model. Results are presented in Figure 2; full results for Qwen are included in Table 4. On MultiBLiMP, the monolingual models generally outperform the multilingual one, except in four Tier 3 languages where the latter shows modest improvements. Compared to Qwen, results are mixed: while it outperforms our multilingual model in most cases, ours remains stronger in eight languages, with no clear trend by tier. On Belebele, both our models perform near chance, while Qwen achieves substantially higher scores across all languages. This pattern extends to most other benchmarks, where Qwen consistently exceeds baseline and outperforms our models on knowledge-intensive and reasoning tasks.

Figure 3 presents results for the bilingual models, focusing on the three best- and worst-performing

tasks from the monolingual models, along with SIB-200. For several tasks—SIB-200, BMLAMA, XCOMPS, and INCLUDE—adding English as a second training language leads to consistent performance gains across most languages. A notable exception is Dutch on INCLUDE, where bilingual training slightly reduces performance. This may be due to domain mismatch: the Dutch corpus includes high-school exam texts, and the addition of English data likely shifts the model away from this domain. Performance on formal linguistic tasks such as MultiBLiMP remains largely unchanged, suggesting that syntactic competence is less sensitive to bilingual input in this setup.

These baseline experiments are intended as a reference point for future work. We envision the BabyBabelLM models and evaluation suite to serve as the basis for a multilingual benchmark challenge in the spirit of BabyLM, encouraging the commu-

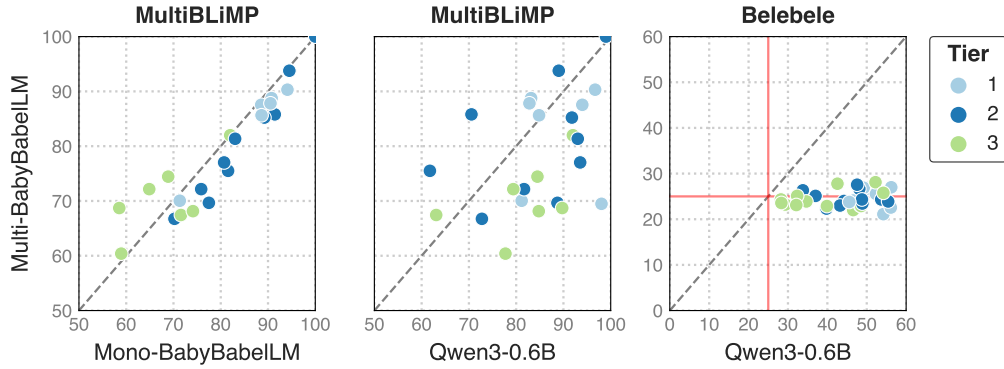


Figure 2: Language-level performance of the multilingual BabyBabelLM model against the monolingual models and Qwen3-0.6B on MultiBLiMP and Belebele. Each point denotes the accuracy on a specific language. Random performance for Belebele is denoted in red.

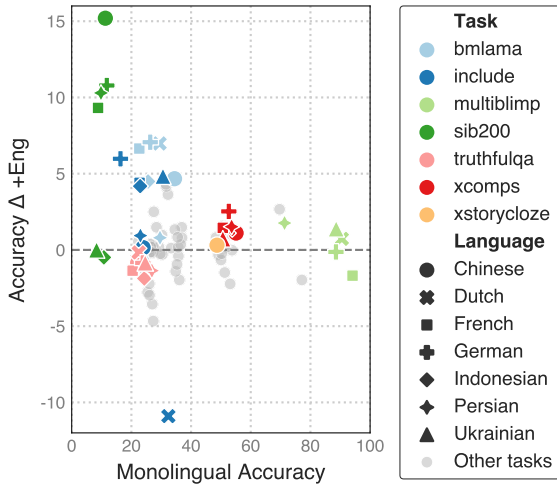


Figure 3: Impact of training LMs on bilingual corpora (adding English) across our evaluation suite. The y-axis denotes the change in accuracy from monolingual to bilingual performance. Dutch SIB-200 performance is omitted for space reasons (+24.8).

nity to explore new training strategies, architectures, and evaluation techniques under comparable, developmentally plausible conditions.

6 Future Outlook

BabyBabelLM is shared research. Starting as a grassroots initiative and conducted in an open and inclusive manner, this resource was gathered by multiple experts with a shared goal. Therefore, we call for this collaboration to continue and welcome further contributions for BabyBabelLM, even after the paper is published.

To act as a complement to the resource, we provide a list of potential questions we believe this resource may aid at answering: Do LMs acquire

language more like language learners of a specific language than another? Are there critical times for learning a second language in LMs (Constantinescu et al., 2025)? Can we replicate results of studies on the border of linguistics and LLMs, where testing on only a single language might bias results (Arnett et al., 2025)? Is there a way to overcome different scripts and unshared tokenizers and provide the same cross-lingual benefits between languages with the same script, or even those that do not share a script? What is the right tokenization scheme across languages, and is tokenization needed at all (Hwang et al., 2025; Rust et al., 2023)? While humans typically give consistent answers across languages, current LMs often do not (Goldman et al., 2025). Even when outputs align, internal changes tend to affect only one language, indicating a degree of separation not seen in human cognition (Ifergan et al., 2024). Can that be changed?

We hope that BabyBabelLM will serve as a foundation for addressing the questions outlined above, and we invite the community to build on this resource to advance a more inclusive and systematic understanding of multilingual language acquisition and modeling.

Limitations

Our resources target a diverse array of audiences, and therefore our decisions are bound to not satisfy each of those perfectly. While deciding between practical constraints, data availability, and potential research needs we prioritized what we believed would make research and experimentation in the BabyLM paradigm easier. Still, we view our dataset only as a starting point. There are many more languages to be included, and even for the

featured languages we imagine further untapped sources of developmentally plausible data.

Despite our language coverage being broader than usual in NLP (cf. Joshi et al., 2020), many languages—particularly those with limited digital presence—remain underrepresented. Especially lacking are languages common in African countries and those with smaller speaker populations, which, despite our efforts, are still underrepresented in our collection. We provide instructions for submitting new languages in our GitHub⁴ and welcome community contributions.

Although we aimed to collect as much cognitively plausible data as possible, we also want to stress that our datasets do not contain the actual language a single native speaker of any of the included languages is exposed to. While our data approximates this input much better than standard pretraining resources (e.g. Wikipedia dumps or datasets like Dolma, Soldaini et al., 2024), the distribution of topics and formats remains only a gross approximation of the diversity experienced by a native learner.

While we calibrate dataset sizes using byte-adjusted thresholds to ensure comparability, the actual composition of developmentally plausible content varies substantially across languages. In several cases, high-quality child-directed speech (CDS) or educational material is unavailable, and we rely more heavily on fallback sources such as subtitles or Wikipedia. This variability may introduce confounds in cross-linguistic analyses and limits the strength of direct typological comparisons. We recommend that future work interpreting model differences across languages take these compositional disparities into account.

Our final limitation is the lack of cross-linguistically available evaluation resources. Many languages are only evaluated on monolingual datasets explicitly created for them, and beyond MultiBLiMP there is currently no resource that covers all included languages. As the study of bilingualism or the acquisition of multiple languages (by models and/or humans) are intended applications of this dataset, we are also constrained by a lack of resources that explicitly target multilingual capabilities. We did not create a standardized testbed to test such questions ourselves, as they are too varied. However, we hope that our data and existing evaluations can serve as inspiration for

further research in that direction.

References

2019. Winogrande: An adversarial winograd schema challenge at scale.
- Ahmed Abdelali, Francisco Guzman, Hassan Sajjad, and Stephan Vogel. 2014. [The AMARA corpus: Building parallel language resources for the educational domain](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, pages 1856–1862, Reykjavik, Iceland. European Language Resources Association (ELRA).
- David Ifeoluwa Adelani, Md Mahfuz Ibn Alam, Antonios Anastasopoulos, Akshita Bhagia, Marta R. Costa-jussà, Jesse Dodge, Fahim Faisal, Christian Federmann, Natalia Fedorova, Francisco Guzmán, Sergey Koshelev, Jean Maillard, Vukosi Marivate, Jonathan Mbuya, Alexandre Mourachko, Safiyyah Saleem, Holger Schwenk, and Guillaume Wenzek. 2022. [Findings of the WMT’22 shared task on large-scale machine translation evaluation for African languages](#). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 773–800, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- David Ifeoluwa Adelani, Hannah Liu, Xiaoyu Shen, Nikita Vassilyev, Jesujoba O. Alabi, Yanke Mao, Haonan Gao, and En-Shiun Annie Lee. 2024. [SIB-200: A simple, inclusive, and big evaluation dataset for topic classification in 200+ languages and dialects](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 226–245, St. Julian’s, Malta. Association for Computational Linguistics.
- Muhammed Al Khalil, Hind Saddiki, Nizar Habash, and Latifa Alfalasi. 2018. [A leveled reading corpus of Modern Standard Arabic](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Latifa Al-Sulaiti, Noorhan Abbas, Claire Brierley, Eric Atwell, and Ayman Alghamdi. 2016. [Compilation of an Arabic children’s corpus](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 1808–1812, Portorož, Slovenia. European Language Resources Association (ELRA).
- Iolanda Alfano, Francesco Cutugno, Aurelio De Rosa, Claudio Iacobini, Renata Savy, Maria Voghera, and 1 others. 2014. Volip: a corpus of spoken italian and a virtuous example of reuse of linguistic resources. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, pages 3897–3901. European Language Resources Association (ELRA).

⁴[link anonymized]

735	Catherine Arnett, Tyler A. Chang, and Benjamin Bergen.	Karissa Vincentio, Jennifer Santoso, David Moel-	789
736	2024. A bit of a problem: Measurement disparities	jadi, Cahya Wirawan, Frederikus Hudi, Muham-	790
737	in dataset sizes across languages . In <i>Proceedings of</i>	mad Satrio Wicaksono, Ivan Parmonangan, Ika Al-	791
738	<i>the 3rd Annual Meeting of the Special Interest Group</i>	fini, Ilham Firdausi Putra, Samsul Rahmadani, and	792
739	<i>on Under-resourced Languages @ LREC-COLING</i>	29 others. 2023a. NusaCrowd: Open source initiative	793
740	2024, pages 1–9, Torino, Italia. ELRA and ICCL.	for Indonesian NLP resources . In <i>Findings of the As-</i>	794
741	Catherine Arnett, Tyler A Chang, James A Michaelov,	<i>sociation for Computational Linguistics: ACL 2023</i> ,	795
742	and Benjamin K Bergen. 2025. On the acquisition	pages 13745–13818, Toronto, Canada. Association	796
743	of shared grammatical representations in bilingual	for Computational Linguistics.	797
744	language models. <i>arXiv preprint arXiv:2503.03962</i> .		
745	Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,	Samuel Cahyawijaya, Holy Lovenia, Fajri Koto, Dea	798
746	Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei	Adhista, Emmanuel Dave, Sarah Oktavianti, Salsabil	799
747	Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin,	Akbar, Jhonson Lee, Nuur Shadieq, Tjeng Wawan	800
748	Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu,	Cenggoro, Hanung Linuwih, Bryan Wilie, Galih	801
749	Keming Lu, and 29 others. 2023. Qwen technical	Muridan, Genta Winata, David Moeljadi, Al-	802
750	report. <i>arXiv preprint arXiv:2309.16609</i> .	ham Fikri Aji, Ayu Purwarianti, and Pascale Fung.	803
751	Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel	2023b. NusaWrites: Constructing high-quality	804
752	Artetxe, Satya Narayan Shukla, Donald Husa, Naman	corpora for underrepresented and extremely low-	805
753	Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and	resource languages . In <i>Proceedings of the 13th In-</i>	806
754	Madian Khabisa. 2024. The belebele benchmark: a	<i>ternational Joint Conference on Natural Language</i>	807
755	parallel reading comprehension dataset in 122 lan-	<i>Processing and the 3rd Conference of the Asia-Pacific</i>	808
756	guage variants . In <i>Proceedings of the 62nd Annual</i>	<i>Chapter of the Association for Computational Lin-</i>	809
757	<i>Meeting of the Association for Computational Lin-</i>	<i>guistics (Volume 1: Long Papers)</i> , pages 921–945,	810
758	<i>guistics (Volume 1: Long Papers)</i> , pages 749–775,	Nusa Dua, Bali. Association for Computational Lin-	811
759	Bangkok, Thailand. Association for Computational	guistics.	812
760	Linguistics.		
761	Ezgi Başar, Francesca Padovani, Jaap Jumelet, and	Thea Cameron-Faulkner, Elena Lieven, and Michael	813
762	Arianna Bisazza. 2025. Turblimp: A turkish	Tomasello. 2003. A construction based analysis of	814
763	benchmark of linguistic minimal pairs . <i>Preprint</i> ,	child directed speech . <i>Cognitive Science</i> , 27(6):843–	815
764	arXiv:2506.13487.	873.	816
765	Yoshua Bengio, Sören Mindermann, and Daniel Privit-		
766	era. 2025. International ai safety report 2025.	Thea Cameron-Faulkner and Claire Noble. 2013. A	817
767	Ryan Bishop. 1991. There’s Nothing Natural About	comparison of book text and Child Directed Speech .	818
768	Natural Conversation: A Look at Dialogue in Fiction	<i>First Language</i> , 33(3):268–279.	819
769	and Drama . <i>Oral Tradition</i> , pages 58–78.		
770	Dominique Brunato. 2025. Learning from impairment:	Luca Capone, Alice Suozzi, Gianluca Lebani, and	820
771	Leveraging insights from clinical linguistics in lan-	Alessandro Lenci. 2024. BaBIes: A benchmark	821
772	guage modelling research . In <i>Proceedings of the 31st</i>	for the linguistic evaluation of Italian baby language	822
773	<i>International Conference on Computational Linguis-</i>	models . In <i>Proceedings of the 10th Italian Confer-</i>	823
774	<i>tics</i> , pages 4167–4174, Abu Dhabi, UAE. Associa-	<i>ence on Computational Linguistics (CLiC-it 2024)</i> ,	824
775	tion for Computational Linguistics.	pages 157–170, Pisa, Italy. CEUR Workshop Pro-	825
776	Bastian Bunzeck and Holger Diessel. 2025. The rich-	ceedings.	826
777	ness of the stimulus: Constructional variation and de-		
778	velopment in child-directed speech . <i>First Language</i> ,	Tyler A. Chang, Catherine Arnett, Zhuowen Tu, and	827
779	45(2):152–176.	Benjamin K. Bergen. 2024. Goldfish: Monolin-	828
780	Bastian Bunzeck, Daniel Duran, and Sina Zarriß. 2025.	gual language models for 350 languages . <i>CoRR</i> ,	829
781	Do construction distributions shape formal language	abs/2408.10441.	830
782	learning in German BabyLMs? In <i>Proceedings of</i>	Lucas Georges Gabriel Charpentier and David Samuel.	831
783	<i>the 29th Conference on Computational Natural Lan-</i>	2023. Not all layers are equally as important: Every	832
784	<i>guage Learning</i> , pages 169–186, Vienna, Austria.	Layer Counts BERT . In <i>Proceedings of the BabyLM</i>	833
785	Association for Computational Linguistics.	<i>Challenge at the 27th Conference on Computational</i>	834
786	Samuel Cahyawijaya, Holy Lovenia, Alham Fikri Aji,	<i>Natural Language Learning</i> , pages 210–224, Singa-	835
787	Genta Winata, Bryan Wilie, Fajri Koto, Rahmad	pore. Association for Computational Linguistics.	836
788	Mahendra, Christian Wibisono, Ade Romadhony,		
		Lucas Georges Gabriel Charpentier and David Samuel.	837
		2024. GPT or BERT: why not both? In <i>The 2nd</i>	838
		<i>BabyLM Challenge at the 28th Conference on Com-</i>	839
		<i>putational Natural Language Learning</i> , pages 262–	840
		283, Miami, FL, USA. Association for Computa-	841
		tional Linguistics.	842
		Jiali Cheng and Hadi Amiri. 2024. Mu-bench: A multi-	843
		task multimodal benchmark for machine unlearning .	844
		<i>Preprint</i> , arXiv:2406.14796.	845

846	Madalina Chitez, Mihai Dascalu, Aura Cristina Udrea,	Mahmoud El-Haj. 2020. Habibi - a multi dialect multi	902
847	Cosmin Strilețchi, Karla Csűrös, Roxana Rogobete,	national Arabic song lyrics corpus . In <i>Proceedings</i>	903
848	and Alexandru Oravițan. 2024. Towards building the	<i>of the Twelfth Language Resources and Evaluation</i>	904
849	LEMI readability platform for children’s literature	<i>Conference</i> , pages 1318–1326, Marseille, France. Eu-	905
850	in the Romanian language . In <i>Proceedings of the</i>	ropean Language Resources Association.	906
851	<i>2024 Joint International Conference on Computa-</i>		
852	<i>tional Linguistics, Language Resources and Evalua-</i>	Ronen Eldan and Yuanzhi Li. 2023. TinyStories: How	907
853	<i>tion (LREC-COLING 2024)</i> , pages 16450–16456,	Small Can Language Models Be and Still Speak Co-	908
854	Torino, Italia. ELRA and ICCL.	herent English? <i>Preprint</i> , arXiv:2305.07759.	909
855	Leshem Choshen, Yang Zhang, and Jacob Andreas.	Steven Y. Feng, Noah Goodman, and Michael Frank.	910
856	2024. A Hitchhiker’s Guide to Scaling Law Esti-	2024. Is child-directed speech effective training	911
857	mation . <i>arXiv preprint</i> .	data for language models? In <i>Proceedings of the</i>	912
858	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,	<i>2024 Conference on Empirical Methods in Natural</i>	913
859	Ashish Sabharwal, Carissa Schoenick, and Oyvind	<i>Language Processing</i> , pages 22055–22071, Miami,	914
860	Tafford. 2018. Think you have solved question	Florida, USA. Association for Computational Lin-	915
861	answering? try arc, the ai2 reasoning challenge .	guistics.	916
862	<i>Preprint</i> , arXiv:1803.05457.	Mathieu Fenniak, Matthew Stamy, pubpub zz, Martin	917
863	Alexis Conneau, Kartikay Khandelwal, Naman Goyal,	Thoma, Matthew Peveler, exiledkingcc, and PyPDF2	918
864	Vishrav Chaudhary, Guillaume Wenzek, Francisco	Contributors. 2022. The PyPDF2 library .	919
865	Guzmán, Edouard Grave, Myle Ott, Luke Zettle-	Achille Fusco, Matilde Barbini, Maria Letizia Pic-	920
866	moyer, and Veselin Stoyanov. 2020. Unsupervised	cini Bianchessi, Veronica Bressan, Sofia Neri, Sarah	921
867	cross-lingual representation learning at scale . In <i>Pro-</i>	Rossi, Tommaso Sgrizzi, and Cristiano Chesi. 2024.	922
868	<i>ceedings of the 58th Annual Meeting of the Asso-</i>	Recurrent networks are (linguistically) better? an	923
869	<i>ciation for Computational Linguistics</i> , pages 8440–	(ongoing) experiment on small-LM training on child-	924
870	8451, Online. Association for Computational Lin-	directed speech in Italian . In <i>Proceedings of the 10th</i>	925
871	guistics.	<i>Italian Conference on Computational Linguistics</i>	926
872	Alexis Conneau, Rutu Rinott, Guillaume Lample, Adina	<i>(CLiC-it 2024)</i> , pages 382–389, Pisa, Italy. CEUR	927
873	Williams, Samuel Bowman, Holger Schwenk, and	Workshop Proceedings.	928
874	Veselin Stoyanov. 2018. XNLI: Evaluating cross-	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Bider-	929
875	lingual sentence representations . In <i>Proceedings of</i>	man, Sid Black, Anthony DiPofi, Charles Foster,	930
876	<i>the 2018 Conference on Empirical Methods in Nat-</i>	Laurence Golding, Jeffrey Hsu, Alain Le Noac’h,	931
877	<i>ural Language Processing</i> , pages 2475–2485, Brus-	Haonan Li, Kyle McDonell, Niklas Muennighoff,	932
878	sels, Belgium. Association for Computational Lin-	Chris Ociepa, Jason Phang, Laria Reynolds, Hailey	933
879	guistics.	Schoelkopf, Aviya Skowron, Lintang Sutawika, and	934
880	Ionut Constantinescu, Tiago Pimentel, Ryan Cotterell,	5 others. 2024. The language model evaluation har-	935
881	and Alex Warstadt. 2025. Investigating critical pe-	ness .	936
882	riod effects in language acquisition through neural	Volker Gast, Christian Wehmeier, and Dirk Vanderbeke.	937
883	language models . <i>Transactions of the Association for</i>	2023. A Register-Based Study of Interior Monologue	938
884	<i>Computational Linguistics</i> , 13:96–120.	in James Joyce’s Ulysses . <i>Literature</i> , 3(1):42–65.	939
885	Yiming Cui, Ting Liu, Zhipeng Chen, Shijin Wang, and	Giuliana Genovese, Maria Spinelli, Leonor J.	940
886	Guoping Hu. 2016. Consensus attention-based neu-	Romero Lauro, Tiziana Aureli, Giulia Castelletti,	941
887	ral networks for Chinese reading comprehension . In	and Mirco Fasolo. 2020. Infant-directed speech as	942
888	<i>Proceedings of COLING 2016, the 26th International</i>	a simplified but not simple register: A longitudinal	943
889	<i>Conference on Computational Linguistics: Technical</i>	study of lexical and syntactic features . <i>Journal of</i>	944
890	<i>Papers</i> , pages 1777–1786, Osaka, Japan. The COL-	<i>Child Language</i> , 47(1):22–44.	945
891	ING 2016 Organizing Committee.	Jill Gilkerson, Jeffrey A. Richards, Steven F. Warren, Ju-	946
892	Yiming Cui, Ting Liu, Ziqing Yang, Zhipeng Chen,	dith K. Montgomery, Charles R. Greenwood, D. Kim-	947
893	Wentao Ma, Wanxiang Che, Shijin Wang, and Guop-	brough Oller, John H. L. Hansen, and Terrance D.	948
894	ing Hu. 2020. A sentence cloze dataset for chinese	Paul. 2017. Mapping the early language environ-	949
895	machine reading comprehension. In <i>Proceedings of</i>	ment using all-day recordings and automated analy-	950
896	<i>the 28th International Conference on Computational</i>	sis . <i>American Journal of Speech-Language Pathol-</i>	951
897	<i>Linguistics (COLING 2020)</i> .	<i>ogy</i> , 26(2):248–265.	952
898	G. William Domhoff and Adam Schneider. 2008. Study-	Omer Goldman, Uri Shaham, Dan Malkin, Sivan Eiger,	953
899	ing dream content using the archive and search engine	Avinatan Hassidim, Yossi Matias, Joshua Maynez,	954
900	on DreamBank.net . <i>Consciousness and Cognition</i> ,	Adi Mayrav Gilady, Jason Riesa, Shruti Rijhwani,	955
901	17(4):1238–1247.	and 1 others. 2025. Eclectic: a novel challenge set	956
		for evaluation of cross-lingual knowledge transfer.	957
		<i>arXiv preprint arXiv:2502.21228</i> .	958

1071	Richard Lastrucci, Jenalea Rajab, Matimba Shingange,	2024. SEACrowd: A multilingual multimodal data	1129
1072	Daniel Njini, and Vukosi Marivate. 2023. Prepar-	hub and benchmark suite for Southeast Asian lan-	1130
1073	ing the vuk’uzenzele and ZA-gov-multilingual South	guages . In <i>Proceedings of the 2024 Conference on</i>	1131
1074	African multilingual corpora . In <i>Proceedings of</i>	<i>Empirical Methods in Natural Language Processing</i> ,	1132
1075	<i>the Fourth workshop on Resources for African In-</i>	pages 5155–5203, Miami, Florida, USA. Association	1133
1076	<i>digenuous Languages (RAIL 2023)</i> , pages 18–25,	for Computational Linguistics.	1134
1077	Dubrovnik, Croatia. Association for Computational		
1078	Linguistics.		
1079	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022a.	Anton Lozhkov, Loubna Ben Allal, Leandro von Werra,	1135
1080	TruthfulQA: Measuring how models mimic human	and Thomas Wolf. 2024. Fineweb-edu: the finest	1136
1081	falsehoods . In <i>Proceedings of the 60th Annual Meet-</i>	collection of educational content .	1137
1082	<i>ing of the Association for Computational Linguistics</i>		
1083	<i>(Volume 1: Long Papers)</i> , pages 3214–3252, Dublin,	Brian MacWhinney. 2000. <i>The CHILDES Project:</i>	1138
1084	Ireland. Association for Computational Linguistics.	<i>Tools for Analyzing Talk</i> , 3 edition. Lawrence Erl-	1139
		baum Associates, Mahwah, NJ.	1140
1085	Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu	Kyle Mahowald, Anna A. Ivanova, Idan A. Blank,	1141
1086	Wang, Shuohui Chen, Daniel Simig, Myle Ott, Na-	Nancy Kanwisher, Joshua B. Tenenbaum, and	1142
1087	man Goyal, Shruti Bhosale, Jingfei Du, Ramakanth	Evelina Fedorenko. 2024. Dissociating language	1143
1088	Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav	and thought in large language models . <i>Trends in</i>	1144
1089	Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettle-	<i>Cognitive Sciences</i> , pages 517–540.	1145
1090	moyer, Zornitsa Kozareva, Mona Diab, and 2 others.		
1091	2022b. Few-shot learning with multilingual gener-	Alexis Matzopoulos, Charl Hendriks, Hishaam Ma-	1146
1092	ative language models . In <i>Proceedings of the 2022</i>	homed, and Francois Meyer. 2025. BabyLMs for	1147
1093	<i>Conference on Empirical Methods in Natural Lan-</i>	isiXhosa: Data-efficient language modelling in a	1148
1094	<i>guage Processing</i> , pages 9019–9052, Abu Dhabi,	low-resource context . In <i>Proceedings of the First</i>	1149
1095	United Arab Emirates. Association for Computa-	<i>Workshop on Language Models for Low-Resource</i>	1150
1096	tional Linguistics.	<i>Languages</i> , pages 240–248, Abu Dhabi, United Arab	1151
		Emirates. Association for Computational Linguistics.	1152
1097	Pierre Lison and Jörg Tiedemann. 2016. OpenSub-	Marina Mayor-Rocher, Cristina Pozo, Nina Melero,	1153
1098	titles2016: Extracting large parallel corpora from	Gonzalo Martínez, María Grandury, and Pedro Re-	1154
1099	movie and TV subtitles . In <i>Proceedings of the Tenth</i>	viriego. 2025. It’s the same but not the same:	1155
1100	<i>International Conference on Language Resources</i>	Do llms distinguish spanish varieties? <i>Preprint</i> ,	1156
1101	<i>and Evaluation (LREC’16)</i> , pages 923–929, Portorož,	arXiv:2504.20049.	1157
1102	Slovenia. European Language Resources Association		
1103	(ELRA).	Cindy McKellar. 2022. Autshumato english-sesotho	1158
1104	Yikang Liu, Yeting Shen, Hongao Zhu, Lilong Xu, Zhi-	parallel corpora . SADIaR Language Resource	1159
1105	heng Qian, Siyuan Song, Kejia Zhang, Jialong Tang,	Repository, License: Creative Commons Attribution	1160
1106	Pei Zhang, Baosong Yang, Rui Wang, and Hai Hu.	4.0 International.	1161
1107	2024. Zhoblmp: a systematic assessment of lan-	Bettina Messmer, Vinko Sabolčec, and Martin Jaggi.	1162
1108	guage models with linguistic minimal pairs in chi-	2025. Enhancing multilingual llm pretraining with	1163
1109	nese . <i>Preprint</i> , arXiv:2411.06096.	model-based data selection . <i>arXiv</i> .	1164
1110	Zhi Liu, Dong Li, Taotao Long, Chaodong Wen, Xian	Francois Meyer and Jan Buys. 2024. Triples-to-	1165
1111	Peng, and Jiaxin Guo. 2025. CSQ: A Chinese Ele-	isiXhosa (T2X): Addressing the challenges of low-	1166
1112	mentary Science Question Dataset with Rich Disci-	resource agglutinative data-to-text generation . In	1167
1113	pline Properties in Adaptive Problem-Solving Pro-	<i>Proceedings of the 2024 Joint International Con-</i>	1168
1114	cess Generation .	<i>ference on Computational Linguistics, Language</i>	1169
1115	Emiddia Longobardi, Clelia Rossi-Arnaud, Pietro	<i>Resources and Evaluation (LREC-COLING 2024)</i> ,	1170
1116	Spataro, Diane L Putnick, and Marc H Bornstein.	pages 16841–16854, Torino, Italia. ELRA and ICCL.	1171
1117	2015. Children’s acquisition of nouns and verbs in	Ludmila Midrigan Ciochina, Victoria Boyd, Lucila	1172
1118	italian: contrasting the roles of frequency and posi-	Sanchez-Ortega, Diana Malanca_Malac, Doina	1173
1119	tional salience in maternal language. <i>Journal of child</i>	Midrigan, and David P. Corina. 2020. Resources	1174
1120	<i>language</i> , 42(1):95–121.	in underrepresented languages: Building a repre-	1175
1121	Holy Lovenia, Rahmad Mahendra, Salsabil Maulana	sentative Romanian corpus . In <i>Proceedings of the</i>	1176
1122	Akbar, Lester James V. Miranda, Jennifer San-	<i>Twelfth Language Resources and Evaluation Confer-</i>	1177
1123	tosso, Elyanah Aco, Akhdan Fadhilah, Jonibek	<i>ence</i> , pages 3291–3296, Marseille, France. European	1178
1124	Mansurov, Joseph Marvin Imperial, Onno P. Kamp-	Language Resources Association.	1179
1125	man, Joel Ruben Antony Moniz, Muhammad	Aaron Mueller, Garrett Nicolai, Panayiota Petrou-	1180
1126	Ravi Shulthan Habibi, Frederikus Hudi, Railey Mon-	Zeniou, Natalia Talmina, and Tal Linzen. 2020.	1181
1127	talan, Ryan Ignatius, Joanito Agili Lopo, William	Cross-linguistic syntactic evaluation of word predic-	1182
1128	Nixon, Börje F. Karlsson, James Jaya, and 42 others.	tion models . In <i>Proceedings of the 58th Annual Meet-</i>	1183
		<i>ing of the Association for Computational Linguistics</i> ,	1184

1185	pages 5523–5539, Online. Association for Computational Linguistics.	1242
1186		1243
1187	Joakim Nivre, Daniel Zeman, Filip Ginter, and Francis Tyers. 2017. Universal Dependencies . In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Tutorial Abstracts</i> , Valencia, Spain. Association for Computational Linguistics.	1244
1188		1245
1189		1246
1190		1247
1191		1248
1192		
1193	Jessica Ojo, Odunayo Ogundepo, Akintunde Oladipo, Kelechi Ogueji, Jimmy Lin, Pontus Stenetorp, and David Ifeoluwa Adelani. 2025. Afrobench: How good are large language models on african languages? <i>Preprint</i> , arXiv:2311.07978.	1249
1194		1250
1195		1251
1196		1252
1197		1253
1198	Francesca Padovani, Jaap Jumelet, Yevgen Matuskevych, and Arianna Bisazza. 2025. Child-directed language does not consistently boost syntax learning in language models . <i>Preprint</i> , arXiv:2505.23689.	1254
1199		1255
1200		1256
1201		1257
1202	Katerina Papantoniou and Yannis Tzitzikas. 2024. Nlp for the greek language: A longer survey . <i>Preprint</i> , arXiv:2408.10962.	1258
1203		
1204		
1205	Guilherme Penedo, Hynek Kydlíček, Vinko Sabolčec, Bettina Messmer, Negar Foroutan, Amir Hossein Kargaran, Colin Raffel, Martin Jaggi, Leandro Von Werra, and Thomas Wolf. 2025. Fineweb2: One pipeline to scale them all – adapting pre-training data processing to every language . <i>Preprint</i> , arXiv:2506.20920.	1259
1206		1260
1207		1261
1208		1262
1209		1263
1210		1264
1211		
1212	Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. 2020. XCOPA: A multilingual dataset for causal common-sense reasoning . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2362–2376, Online. Association for Computational Linguistics.	1265
1213		1266
1214		1267
1215		1268
1216		1269
1217		
1218		
1219	Velka Popova. 2020. Childes bulgarian labling corpus .	1270
1220		1271
1221	Laurent Prévot, Sheng-Fu Wang, Jou-An Chi, and Shu-Kai Hsieh. 2024. Extending the BabyLM initiative : Promoting diversity in datasets and metrics through high-quality linguistic corpora . In <i>The 2nd BabyLM Challenge at the 28th Conference on Computational Natural Language Learning</i> , pages 147–158, Miami, FL, USA. Association for Computational Linguistics.	1272
1222		1273
1223		1274
1224		1275
1225		1276
1226		1277
1227	Ayu Purwarianti, Dea Adhista, Agung Baptiso, Miftahul Mahfuzh, Yusrina Sabila, Aulia Adila, Samuel Cahyawijaya, and Alham Fikri Aji. 2025. NusaDialogue: Dialogue summarization and generation for underrepresented and extremely low-resource languages . In <i>Proceedings of the Second Workshop in South East Asian Language Processing</i> , pages 82–100, Online. Association for Computational Linguistics.	1278
1228		1279
1229		1280
1230		1281
1231		1282
1232		1283
1233		1284
1234		1285
1235		1286
1236	Jirui Qi, Raquel Fernández, and Arianna Bisazza. 2023. Cross-lingual consistency of factual knowledge in multilingual language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 10650–10666, Singapore. Association for Computational Linguistics.	1287
1237		1288
1238		1289
1239		1290
1240		1291
1241		1292
		1293
		1294
		1295
		1296
		1297
		1298
		1299
	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, and 25 others. 2025. Qwen2.5 technical report . <i>Preprint</i> , arXiv:2412.15115.	
	Angelika Romanou, Negar Foroutan, Anna Sotnikova, Zeming Chen, Sree Harsha Nelaturu, Shivalika Singh, Rishabh Maheshwary, Micol Altomare, Mohamed A. Haggag, Snegha A, Alfonso Amayuelas, Azril Hafizi Amirudin, Viraat Aryabumi, Danylo Boiko, Michael Chang, Jenny Chim, Gal Cohen, Aditya Kumar Dalmia, Abraham Diress, and 40 others. 2024. Include: Evaluating multilingual language understanding with regional knowledge . <i>Preprint</i> , arXiv:2411.19799.	
	Dimitris Roussis, Leon Voukoutis, Georgios Paraskevopoulos, Sokratis Sofianopoulos, Prokopis Prokopidis, Vassilis Papavasileiou, Athanasios Katsamanis, Stelios Piperidis, and Vassilis Katsourous. 2025. Krikri: Advancing open large language models for greek . <i>Preprint</i> , arXiv:2505.13772.	
	Phillip Rust, Jonas F. Lotz, Emanuele Bugliarello, Elizabeth Salesky, Miryam de Lhoneux, and Desmond Elliott. 2023. Language modelling with pixels . In <i>The Eleventh International Conference on Learning Representations</i> .	
	Suchir Salhan, Richard Diehl Martinez, Zébulon Goriely, and Paula Buttery. 2024. Less is more: Pre-training cross-lingual small-scale language models with cognitively-plausible curriculum learning strategies . In <i>The 2nd BabyLM Challenge at the 28th Conference on Computational Natural Language Learning</i> , pages 174–188, Miami, FL, USA. Association for Computational Linguistics.	
	Guokan Shang, Hadi Abdine, Yousef Khoubrane, Amr Mohamed, Yassine Abbahaddou, Sofiane Ennadir, Imane Momayiz, Xuguang Ren, Eric Moulines, Preslav Nakov, Michalis Vazirgiannis, and Eric Xing. 2025. Atlas-chat: Adapting large language models for low-resource Moroccan Arabic dialect . In <i>Proceedings of the First Workshop on Language Models for Low-Resource Languages</i> , pages 9–30, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
	Matthew Shardlow, Fernando Alva-Manchego, Riza Batista-Navarro, Stefan Bott, Saul Calderon Ramirez, Rémi Cardon, Thomas François, Akio Hayakawa, Andrea Horbach, Anna Huelsing, Yusuke Ide, Joseph Marvin Imperial, Adam Nohejl, Kai North, Laura Occhipinti, Nelson Pérez Rojas, Nishat Raihan, Tharindu Ranasinghe, Martin Solis Salazar, and 2 others. 2024. An Extensible Massively Multilingual Lexical Simplification Pipeline Dataset using the MultiLS Framework . In <i>Proceedings of the 3rd Workshop on Tools and Resources for People with READING Difficulties (READI)</i> .	

1300	Zhewen Shen, Aditya Joshi, and Ruey-Cheng Chen.	Maria Spinelli, Chiara Suttora, Adrian Garcia-Sierra,	1355
1301	2024. BAMBINO-LM: (bilingual-)human-inspired	Fabia Franco, Francesca Lionetti, and Mirco Fasolo.	1356
1302	continual pre-training of BabyLM . In <i>Proceedings</i>	2023. Editorial: Are there different types of child-	1357
1303	<i>of the Workshop on Cognitive Modeling and Compu-</i>	directed speech? dynamic variations according to	1358
1304	<i>tational Linguistics</i> , pages 1–7, Bangkok, Thailand.	individual and contextual factors . <i>Frontiers in Psy-</i>	1359
1305	Association for Computational Linguistics.	<i>chology</i> , Volume 13 - 2022.	1360
1306	Maria Shvedova and Arsenii Lukashevskiy. 2024. Plug:	Alice Suozzi, Luca Capone, Gianluca E Lebani, and	1361
1307	Corpus of old ukrainian texts . Available at https://github.com/Dandelliony/pluperfect_grac .	Alessandro Lenci. 2025. Bambi: Developing	1362
1308		baby language models for italian . <i>arXiv preprint</i>	1363
1309	Johannes Sibeko and Menno Zaanen. 2023. A data set	<i>arXiv:2503.09481</i> .	1364
1310	of final year high school examination texts of south	Shira Tal, Eitan Grossman, and Inbal Arnon. 2024.	1365
1311	african home and first additional language subjects .	Infant-directed speech becomes less redundant as in-	1366
1312	<i>Journal of Open Humanities Data</i> , 9.	fants grow: Implications for language learning . <i>Cog-</i>	1367
1313	Mariana O Silva, Clarisse Scofield, and Mirella M Moro.	<i>nition</i> , 249:105817.	1368
1314	2021. Pportal: Public domain portuguese-language	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya	1369
1315	literature dataset . In <i>Dataset Showcase Workshop</i>	Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin,	1370
1316	<i>(DSW)</i> , pages 77–88. SBC.	Tatiana Matejovicova, Alexandre Ramé, Morgane	1371
1317	Shivalika Singh, Angelika Romanou, Clémentine Four-	Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey	1372
1318	rier, David Ifeoluwa Adelani, Jian Gang Ngui, Daniel	Cideron, Jean bastien Grill, Sabela Ramos, Edouard	1373
1319	Vila-Suero, Peerat Limkonchotiawat, Kelly Marchi-	Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev,	1374
1320	sio, Wei Qi Leong, Yosephine Susanto, Raymond	and 197 others. 2025. Gemma 3 technical report .	1375
1321	Ng, Shayne Longpre, Sebastian Ruder, Wei-Yin	<i>Preprint</i> , arXiv:2503.19786.	1376
1322	Ko, Antoine Bosselut, Alice Oh, Andre Martins,	Agnes Tellings, Micha Hulsbosch, Anne Vermeer, and	1377
1323	Leshem Choshen, Daphne Ippolito, and 4 others.	Antal van den Bosch. 2014. Basilex: an 11.5 mil-	1378
1324	2025. Global MMLU: Understanding and addressing	lion words corpus of dutch texts written for children .	1379
1325	cultural and linguistic biases in multilingual evalua-	<i>Computational Linguistics in the Netherlands Jour-</i>	1380
1326	tion . In <i>Proceedings of the 63rd Annual Meeting of</i>	<i>nal</i> , 4:191–208.	1381
1327	<i>the Association for Computational Linguistics (Vol-</i>	Jannis Vamvas and Rico Sennrich. 2021. On the lim-	1382
1328	<i>ume 1: Long Papers)</i> , pages 18761–18799, Vienna,	its of minimal pairs in contrastive evaluation . In	1383
1329	Austria. Association for Computational Linguistics.	<i>Proceedings of the Fourth BlackboxNLP Workshop</i>	1384
1330	Dan Isaac Slobin. 2014. <i>The crosslinguistic study of</i>	<i>on Analyzing and Interpreting Neural Networks for</i>	1385
1331	<i>language acquisition: Volume 5: Expanding the con-</i>	<i>NLP</i> , pages 58–68, Punta Cana, Dominican Republic.	1386
1332	<i>texts</i> . Psychology Press.	Association for Computational Linguistics.	1387
1333	R. Smith. 2007. An overview of the tesseract ocr engine .	Leon Voukoutis, Dimitris Roussis, Georgios	1388
1334	In <i>Ninth International Conference on Document Anal-</i>	Paraskevopoulos, Sokratis Sofianopoulos, Prokopis	1389
1335	<i>ysis and Recognition (ICDAR 2007)</i> , volume 2, pages	Prokopidis, Vassilis Papavasileiou, Athanasios	1390
1336	629–633.	Katsamanis, Stelios Piperidis, and Vassilis Katsouras.	1391
1337	Catherine E. Snow and Charles A. Ferguson, editors.	2024. Meltemi: The first open large language model	1392
1338	1977. <i>Talking to Children: Language Input and Ac-</i>	for greek . <i>Preprint</i> , arXiv:2407.20743.	1393
1339	<i>quisition</i> . Cambridge University Press, Cambridge,	Xiaoyang Wang, Chen Li, Jianqiao Zhao, and Dong Yu.	1394
1340	MA.	2021. Naturalconv: A chinese dialogue dataset to-	1395
1341	Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin	wards multi-turn topic-driven conversation . <i>Proceed-</i>	1396
1342	Schwenk, David Atkinson, Russell Authur, Ben Bo-	<i>ings of the AAAI Conference on Artificial Intelligence</i> ,	1397
1343	gin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar,	35(16):14006–14014.	1398
1344	Valentin Hofmann, Ananya Harsh Jha, Sachin Kumar,	Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan	1399
1345	Li Lucy, Xinxin Lyu, Nathan Lambert, Ian Magnusson,	Wilcox, Chengxu Zhuang, Juan Ciro, Rafael Mos-	1400
1346	Jacob Morrison, Niklas Muennighoff, and 17 others.	quera, Bhargavi Paranjabe, Adina Williams, Tal	1401
1347	2024. Dolma: An Open Corpus of Three Trillion	Lizen, and Ryan Cotterell. 2023. Findings of the	1402
1348	Tokens for Language Model Pretraining Research .	BabyLM challenge: Sample-efficient pretraining on	1403
1349	<i>Preprint</i> , arXiv:2402.00159.	developmentally plausible corpora . In <i>Proceedings</i>	1404
1350	Taiga Someya and Yohei Oseki. 2023. JBLiMP:	<i>of the BabyLM Challenge at the 27th Conference on</i>	1405
1351	Japanese benchmark of linguistic minimal pairs . In	<i>Computational Natural Language Learning</i> , pages	1406
1352	<i>Findings of the Association for Computational Lin-</i>	1–34, Singapore. Association for Computational Lin-	1407
1353	<i>guistics: EACL 2023</i> , pages 1581–1594, Dubrovnik,	guistics.	1408
1354	Croatia. Association for Computational Linguistics.	Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mo-	1409
		hananey, Wei Peng, Sheng-Fu Wang, and Samuel R.	1410

dialects. Even though spoken Arabic can vary substantially across different regions, with speech being often mutually unintelligible, due to the scarcity of developmentally plausible data we opted to combine all the different dialects into one dataset.

In recent years, there has been substantial effort towards advancing NLP for the Arabic languages. However, child-oriented resources and developmentally plausible corpora are still lacking. Some efforts have been made such as a Compilation of an Arabic Children’s Corpus (Al-Sulaiti et al., 2016), and a leveled reading corpus of Modern Standard Arabic (Al Khalil et al., 2018). However, the data has not been made publicly available. Additionally, natural conversation and speech data, which is part of the child’s linguistic environment, is often unreleased or available only under a fee.

Despite these restrictions, we present a developmentally plausible dataset for Arabic, consisting of children’s books and stories, song lyrics, natural conversations, and articles from child wikis. For books, we include the recently released Arabic Book Corpus (Hallberg, 2025), keeping only the “children stories” category, containing both translated and original titles, mostly from the 20th century. Given that songs form a common linguistic input for children, we incorporate in our data the Habibi Corpus (El-Haj, 2020), consisting of song lyrics in a variety of dialects. Adult speech data was collected from MagicHub⁵ for the Yemeni⁶ and the Egyptian⁷ dialects, as natural conversation is present in the child’s learning environment. Lastly, we include children stories from GlotStoryBooks and Ririro, as well child directed speech from CHILDES, and Children’s Wiki articles.

Given the large variety and differences between Arabic dialects, our future goal is to collect more dialect-specific resources, and release single dialect developmentally plausible datasets, incorporating data from recent efforts such as Atlas-Chat (Shang et al., 2025). This will give us the opportunity to study the unique nature of dialects in the Arabic language, and how they interact in terms of language model training and performance in a developmentally plausible setting.

⁵<https://magichub.com/datasets/>

⁶Available here: <https://magichub.com/datasets/yemeni-arabic-conversational-speech-corpus/>

⁷Available here: <https://magichub.com/datasets/egyptian-arabic-conversational-speech-corpus/>

B.2 Bulgarian

The Bulgarian dataset is a compilation of children’s literature accessed via a public library website: <https://chitanka.info>. The Bulgarian BabyLM corpus is the first large-scale corpus of child-appropriate Bulgarian text.

Dataset Description. The Chitanka portion of the Bulgarian BabyLM corpus consists of 28M tokens, excluding punctuation. To our knowledge, the only other similarly sourced dataset is from CHILDES, which is also included in the Bulgarian BabyLM Corpus. The Bulgarian portion of CHILDES consists of 94K tokens of Child-Directed Speech (CDS) collected by Popova (2020) for 5 children aged 1-2 years. The Chitanka library consists of several categories of books, ranging from science to literature, and has a curated section of Children and Young Adult’s literature that the site owners has confirmed are free to distribute.⁸ Chitanka’s Children and Young Adult Literature sections consists of 670 texts, comprised of novels and short stories (377 texts) poems and riddles (40 texts); fairy tales (169 texts); and other children stories (25 texts); and miscellaneous children and young adult literature (68 texts).⁹

Preprocessing. The individual texts have been cleaned of front and back matter. Each text is provided alongside the link it was scraped from. The largest version of the dataset includes children’s literature for various ages, but if one would like to restrict the dataset to a subset of earlier-age appropriate literature, this can be done by restricting the URLs which correspond to the Bulgarian Ministry of Education’s programme for second and third grade summer reading (ages 6-8 years), for which the corresponding URLs are listed in the README of the dataset. A final notable detail of the dataset is that the texts are in the Cyrillic alphabet, which should be considered during preprocessing.

B.3 Cantonese

We compile our Cantonese text corpus by consolidating four publicly available resources: the Hambaanglaang project, the GlotStory Book project,

⁸Excerpt from author correspondence: “Everything in my library is completely free; you’re welcome to use any of the available resources. The books we add are supposed to be free of copyright claims. If such claims do arise—that is, if rights holders or distributors get in touch with us—those works are ‘quarantined’ until a previously agreed period of time has elapsed.”

⁹Available here: <https://chitanka.info/books/category/detska-literatura>

and two Cantonese datasets from CHILDES (HKU-70 Corpus and Lee/Wong/Leung Corpus).

Hambaanglaang¹⁰ An open-source repository of Cantonese graded readers created by volunteers. It offers a collection of stories designed for children across five proficiency levels, aiming to support Cantonese literacy and reading skills within the community. Detailed information about this project can be found in its official documentation.

GlottStory Book (Hong Kong edition)¹¹ A free, open-source literacy site that localizes 40 children’s stories—originally from the African Storybook Project—into multiple languages used in Hong Kong’s “two scripts / three languages” environment (spoken & written Cantonese, Mandarin, English). Here we only extracted Cantonese. Each story is tagged with one of five length/lexical-complexity levels and accompanied by narrated audio recordings intended to support family-, school-, and community-based language learning.

HKU-70 Corpus¹² This dataset contains 70 transcripts of interviews with 70 Cantonese-speaking children aged 2 years 6 months to 5 years 6 months. The data were collected at the University of Hong Kong and represent naturalistic child–adult interactions in preschool settings. Each child participated in a one-hour recording session, with conversations organized around familiar daily routines (e.g., bathing, dressing, feeding) to elicit a diverse range of utterances and syntactic structures. The sample was balanced by gender, and all children were prescreened using the Cantonese version of the Reynell Developmental Language Scales.

Lee/Wong/Leung Corpus¹³ This dataset provides longitudinal data on eight Cantonese-speaking children, each recorded for approximately one year. The recordings capture natural interactions between children, their caregivers, and occasionally other adults. Detailed metadata about the children, including their ages and family backgrounds, are included, providing valuable sociolinguistic context for the dataset.

Preprocessing All datasets were cleaned to retain only complete Traditional Chinese text, with non-textual annotations such as speaker labels and syntactic tags removed. Each storybook page

or conversational block is treated as a single passage-level entry. All text is tokenized using the Qwen1.5-7B (Bai et al., 2023) tokenizer to ensure compatibility with downstream language modeling tasks.

B.4 Mandarin Chinese

In addition to multilingual resources (CHILDES and GlottStoryBook), our Mandarin Chinese dataset includes data from multiple resources. We used data from wikipedia for padding.

Children Speech We first incorporate transcriptions in ChildMandarin (Zhou et al., 2024). This dataset contains high-quality speech data collected from 397 children in China, along with carefully crafted, character-level manual transcriptions.

Children’s Book We use children’s book and stories data collected from various sources. We first obtained children’s stories from Quangushi (full stories)¹⁴. Then, we used children’s stories from and two Chinese reading comprehension datasets: CFT (Cui et al., 2016) and CMRC-2019 (Cui et al., 2020). These two datasets are respectively Cloze and sentence ordering benchmarks derived from children’s stories. We reconstructed the complete stories using the answers provided by the authors and included the stories in our dataset. We also collected open-source children’s book and children’s wiki data from WikiJunior and Wikibooks.

Educational Materials For educational materials, we used several datasets that evaluate models’ general knowledge through exam-style questions, as these datasets are typically well-documented and come with openly available licenses:

- **GAOKAO (Zhang et al., 2023)**: an evaluation framework that uses Chinese National College Entrance Examination (GAOKAO) questions as a dataset to evaluate LLMs. The dataset includes subjective and objective questions from exams from 2010 to 2024.
- **CK-12 (You et al., 2024)**: an evaluation for Chinese LLMs. constructed based on multi-level knowledge graph and covers most comprehensive knowledge points in Chinese K12 field.
- **CSQ (Liu et al., 2025)**: a Chinese Science Question dataset covering four subjects and multiple topics at the Chinese primary school.

¹⁰<https://hambaanglaang.hk/about-us-2/>

¹¹<https://global-asp.github.io/storybooks-hongkong>

¹²<https://talkbank.org/childes/access/Chinese/Cantonese/HKU.html>

¹³<https://talkbank.org/childes/access/Chinese/Cantonese/LeeWongLeung.html>

¹⁴<http://quangushi.com/>

We included the full question prompts, answer choices, correct answers, and explanations. Questions in English are excluded. In addition, we collected grammatical and corrected sentences from FCGEC (Xu et al., 2022), a human-annotated corpus based on multi-choice grammatical error problems. We also collected data from a hierarchical corpus of primary school students’ compositions (Zhou and Zheng, 2024). The primary source of this corpus is elementary school student composition magazines, which ensures that the essays are of relatively high quality.

Conversation Dataset NaturalConv (Wang et al., 2021), a multi-turn topic-driven conversation dataset is also included as such conversations on daily topic are considered as children available.

B.5 Dutch

The Dutch data is built from various educational sources. Licensing laws are very strictly defined in the Netherlands, which makes it challenging to find children’s literature with creative commons licenses. Educational resources, however, are often released under CC-BY license. We include the texts of all high school exams¹⁵ from 1999 to 2024, for all Dutch high school levels: VMBO (age 15–16), HAVO (age 16–17), and VWO (age 17–18), resulting in 6.87M tokens. We extracted 8.78M tokens from WikiWijs¹⁶, a platform for sharing educational materials by teachers for both primary and high school level. KlasCement¹⁷ provides a similar platform, focused on Flemish education, from which we extract 0.14M tokens. Next to these educational resources, we also incorporate BasiLex (Tellings et al., 2014) into the Dutch section. BasiLex contains a collection of child-directed resources, extracted from children’s media, children’s books, and educational materials. We collect 11.37M tokens from BasiLex.

B.6 French

In addition to child-directed speech from CHILDES (around 4 million tokens), we include the following developmentally plausible resources, covering a range of spoken and written language that children are likely to hear or read:

Subtitles (around 6 million tokens): This portion is made mostly of subtitles from the popular animated series Caillou, aimed at toddlers and shared on YouTube by the channel Caillou Français – WildBrain.¹⁸ It includes 1,539 video episodes. In addition to Caillou, we included other well-known children’s shows in France, such as Olive et Tom (171 videos), Lou (52 videos), La vie (8 videos), and a few other youth-oriented clips (15 videos). Each subtitle document includes the YouTube video ID so the original video can be accessed. We obtained raw transcripts via the YouTubeTranscriptApi,¹⁹ filtered for manually entered transcripts (as opposed to automatically created ones), and fed them through the library’s built-in TextFormatter, which strips out all timing information and reassembles each subtitle fragment as plain text.

Conversational data (around 2 million tokens): We include transcripts of spoken conversations that are not directly addressed to children, but that children could realistically overhear. We selected a number of sources from Claire-Dialogue-French-0.1²⁰ (Hunter et al., 2023), including three types of settings: spontaneous everyday conversations (in homes, cafés, or on the street), guided one-on-one interviews, and workplace meetings. The data comes from sources like PFC_free, OFROM, CLAPI, ORFEO_coralrom, ParisStories, CFPP, ACSYNT, and ORFEO_reunions_de_travail.

Children’s books (around 1 million tokens): This portion includes eighteen children’s books.²¹ These were selected to match the reading level of children aged 6 to 12 and to cover a variety of story types. The collection includes classic fairy tales (Contes de Perrault, Grimm, Andersen), simple educational texts (Abécédaire du petit naturaliste, Histoires comme ça pour les petits), and famous adventure and fantasy stories like Le Tour du monde en quatre-vingts jours, L’Île au trésor, Alice au pays des merveilles, and Croc-Blanc.

¹⁵Released publicly by examenblad.nl, with archives available at alleexamens.nl.

¹⁶wikiwijs.nl

¹⁷klascement.net

¹⁸<https://www.youtube.com/@CaillouFrench>

¹⁹<https://pypi.org/project/youtube-transcript-api/>

²⁰<https://huggingface.co/datasets/OpenLLM-France/Claire-Dialogue-French-0.1>

²¹Hand-picked from the Wikisource category: [Catégorie:Littérature jeunesse](#)

B.7 German

Our German data builds upon the existing German BabyLM corpus by [Bunzeck et al. \(2025\)](#), extending it with more developmentally plausible data and discarding the majority of their padding data in favor of the multilingual padding data compiled in the current project. As German is a comparatively high-resource language, we are able to supplement the multilingual resources with a variety of monolingual corpora. Five different children’s wikis are available for German, including the state-sponsored Klexikon and MiniKlexikon, but also comparable efforts for Austrian German like the Kiwithek. In addition, we supplement this kind of educational data with the WikiBooks Wikijunior bookshelf, which is fairly comprehensive for German.

As for books, we aim to make an educated selection of the Project Gutenberg collection featuring works for children and young adults. We include books that are considered classics of children’s literature and read to this day. We further also include classics of German literature that are regularly read in middle school (e.g. works by Franz Kafka). Although they are located at the end of the ‘developmentally plausible’ timeline, they are plausibly encountered by many young adults in the German education system. Similarly, we also include the archives of the *Fluter* magazine, published by the German Federal Agency for Civic Education, which contains a large body of non-fiction writing aimed at adolescents and young adults.

Moving from child-directed to child-available language, we furthermore incorporate the German section of the CallHome corpus ([Karins et al., 1997](#)), which contains transcribed telephone conversations between adults. Such conversations could i) be plausibly overheard by children, and ii) approximate child-directed input nicely by being transcribed from spoken data, which differs quite dramatically from written data in composition (cf. [Cameron-Faulkner et al., 2003](#); [Cameron-Faulkner and Noble, 2013](#); [Bunzeck and Diessel, 2025](#)). In a similar vein, we also incorporate the German portion of Dreambank ([Domhoff and Schneider, 2008](#)), a large corpus of dream reports by adults and children. Despite not being originally spoken, the ‘self-reporting’ register included in this data is closer to spoken data than ordinary writing, and social storytelling is an important component of language acquisition. Therefore, we conclude that

this dataset also enhances the variety and developmental plausibility of the German data.

B.8 Greek

NLP for the Greek language has developed drastically over the past few years ([Papantoniou and Tzitzikas, 2024](#)), with a notable example being the recent release of large language models for the Greek language: Meltemi 7B ([Voukoutis et al., 2024](#)) the first such open LLM for Greek, and Krikri 8B ([Roussis et al., 2025](#)) further scaling up data and model sizes. Here we present, to our knowledge, the first developmentally plausible corpus for the Greek language. The data is curated as a collection of publicly available datasets, sourced mostly from CLARIN:EL²², and original web-scraped children’s books and stories. We present details about the dataset composition and preprocessing below. In the future we plan to include more child-directed speech data in collaboration with language acquisition researchers, incorporating efforts such as the Greek Children Spoken Language Corpus²³ and the Greek-speaking Children Corpus²⁴.

Educational Corpora We incorporate into our data a variety of educational textbooks. We include a selection of Primary School Books²⁵ in the fields of arts, language, religion, history, and social and political sciences, aimed at grades 1-6 (ages 6-12). Apart from textbooks aimed at children, we decided to additionally include material designed for later grades and ages. Even though this content is aimed at the tail end of our target ages for developmentally plausible corpora, we consider it sufficiently relevant and representative of the linguistic input of children. Thus, we collect the CGL Modern Greek Texts corpus²⁶, which comprises around 2 million words from textbooks published by the Greek Ministry of Education taught through grades 7-12 (ages 13-18) in the public school system. We also include the corpus of Pedagogical Greek L2 textbooks²⁷, addressed to indigenous populations or minorities learning Greek as a second language, aimed at proficiency levels A1 to C2 and ages 6-18+. Even though this resource is designed for

²²<https://inventory.clarin.gr/>

²³<http://gcs1.ece.uth.gr/>

²⁴<https://gavriilidou.gr/greek-speaking-children-corpus/>

²⁵<https://inventory.clarin.gr/corpus/1075>

²⁶<http://hdl.handle.net/11500/>

KEG-0000-0000-24FD-B

²⁷<http://hdl.handle.net/11500/ATHENA-0000-0000-2631-E>

non-native learners, we believe the material to be sufficiently close in nature to the learning resources for native Greek speakers. Finally, we include articles from Children’s Wikis.

Children Books Numerous websites host open access e-books and children stories for the Greek language. We identified [openbook.gr](https://www.openbook.gr)²⁸ and free-ebooks.gr²⁹ as the largest such sites, and manually scraped them, selecting e-books from the categories of children, young-adult, and preschool-education. The data consists of children books in the Public Domain, as well as open access books released with permissive licenses. We also include a collection of children stories scraped from [paidika-paramythia.gr](https://www.paidika-paramythia.gr)³⁰. The site enables any author to make a submission in collaboration with the moderators, and includes stories from tradition and mythology, as well as original entries. Lastly, we include sort stories provided in the GlotStoryBooks corpus.

Child Speech We collect publicly available data corresponding to child-produced and child-directed speech. Child Speech³¹ contains transcriptions of children’s speech with a focus on narration; as the result of interviews conducted by university students with children related to them either by friendship or kinship. Our second addition is the Greek Student Chat Dataset³² consists of chat among students (grades 4-18) in online collaborative learning environments (wikis). Finally, we also include the Greek portion of CHILDES noting that speech is transcribed in the Latin script. In future efforts we plan on either removing this data or transliterating it to the Greek script.

Adult Conversation Data Everyday conversations between adults are a natural stimulus for children during language development. We include in our data a corpus of written transcripts of everyday conversations between students of the Department of Linguistics³³ that took place between 2001 and 2006. The data is further supplemented by the

Babinotis archive³⁴, consisting of the same data variety recorded in 2020. The speech is authentic and idiomatic with speakers labeled, resulting in a high quality spoken Greek corpus.

Preprocessing. Regarding e-books, processing the text proved challenging, and required a substantial amount of manual labor. Initially licensing information was extracted, and the corpus was filtered to include only permissive licenses (e.g., cc-by-nc). For [openbook.gr](https://www.openbook.gr), license information is provided as metadata for each entry, while for free-ebooks.gr we manually annotate each book with its license as stated in the text. The stories in [paidika-paramythia.gr](https://www.paidika-paramythia.gr) are released as Public Domain. The text is first extracted from e-books using PyMuPDF³⁵, and is then filtered to remove license statements, author biographies, and other information deemed irrelevant. Further document-specific normalization follows, fixing text extraction errors, removing unwanted unicode characters, and ensuring the validity of the book content. As part of this process, documents deemed unsuitable for children are excluded. For the public datasets in our corpus, standard pre-processing was applied. Morphological and other linguistic annotations were removed from speech data. We note that to ensure anonymity, placeholders exist in conversational text that substitute real information (e.g., names, locations). Lastly, the Primary School Books corpus required considerable cleaning and normalization efforts, containing web-scraping artifacts such as javascript code.

B.9 Italian

In addition to multilingual resources, we have access to transcripts from studies on child language acquisition (Longobardi et al., 2015; Whittle and Nuzzo, 2015; Spinelli et al., 2023), which have already been used as training data in the work of (Suozzi et al., 2025). Regarding book selection, we are able to include approximately thirty books from the independent Italian publishing house Biancoenero Edizioni³⁶, which kindly shares them with us upon request. This publisher has long been committed to the Alta Leggibilità (“High Readability”) project, aimed at making books accessible to all children, including those with reading difficulties. All books are written by Italian authors

²⁸<https://www.openbook.gr/literature>

²⁹<https://free-ebooks.gr/tag/16?>

³⁰<https://www.paidika-paramythia.gr/16>

³¹<http://hdl.handle.net/11500/CLARIN-EL-0000-0000-610D-5>

³²<http://hdl.handle.net/11500/IONION-0000-0000-5E14-1>

³³<http://hdl.handle.net/11500/UOA-0000-0000-5D9C-9>

³⁴<http://hdl.handle.net/11500/UOA-0000-0000-2515-F>

³⁵<https://github.com/pymupdf/PyMuPDF>

³⁶<https://www.biancoeneroedizioni.it/>

and are targeted at readers between the ages of 4 and 10. The themes span a range of topics including environment and ecology, bullying, mystery, diversity and inclusion, growing up and intergenerational relationships, and adventure, according to the categories listed in the publisher’s updated catalog. We also incorporate books from the Logos Group library³⁷. This collection comprises classic children’s stories and fairy tales authored by both Italian and foreign writers whose works are translated into Italian. The estimated target reading age for these texts ranges from approximately 6 to 14 years; however, some of these stories may be orally presented to younger children. Furthermore, we include a series of fairy tales (all from copyright expired sources) curated by the researchers in this study (Fusco et al., 2024). The book section concludes with a manually curated selection of approximately 50 titles from the Project Gutenberg catalog. These works are either explicitly included in the national curriculum for lower and upper secondary education, or authored by canonical figures whose writings are frequently excerpted in educational contexts and whose titles are broadly recognized within the Italian school system. These include both Italian and non-Italian authors. Although the language used in these works is occasionally archaic and stylistically distant from contemporary Italian — as similarly observed in the case of German (and other languages, where applicable) — their inclusion aligns with the upper boundary of the "developmentally plausible" timeline. Nevertheless, these texts remain realistically encountered by a substantial portion of young adults within the Italian educational system. Leveraging a dataset previously curated by (Fusco et al., 2024), we incorporate a collection of children’s songs from the Zecchino D’Oro archive, a renowned and long-standing Italian music festival for children.

As for educational resources, our dataset includes the Italian portion of the WikiBooks Wikijunior bookshelf, which comprises a variety of entries covering a diverse set of topics (e.g., the human body, dinosaurs, the solar system). This section also includes around 60 YouTube video transcripts from the animated cartoon Calimero³⁸. This subset of videos is selected based on the presence of consistent and realistic punctuation, while all automatically generated transcripts containing grammatical

errors and typos are filtered out. The educational materials section also includes a catalog of past INVALSI assessments of Italian and Math at both primary and secondary school levels³⁹. INVALSI is the national body responsible for evaluating student competencies and the quality of the Italian education system. Lastly, we include an archive of national high school final examination prompts made available by the Italian Ministry of Education⁴⁰, covering the past 20 years. These standardized exams, taken by all students aged 18–19 to obtain their diploma, vary by school type. Although situated at the final stage of secondary education, we consider this material relevant, as it reflects the curricular exposure of the vast majority of Italian students and offers a representative snapshot of the competencies expected of young adults within the national education system.

As supplementary resources, we include two text simplification datasets. The first, from (Brunato, 2025), comprises Terence and Teacher: Terence contains 32 short Italian children’s stories with expert-produced simplifications for readers with comprehension difficulties, while Teacher includes 18 pairs of texts from various genres (e.g., literature, textbooks) used in educational settings. The second dataset, MultiLS, is developed for the MLSP2024 shared task by (Shardlow et al., 2024).

Finally, with respect to adult spoken Italian that remains accessible to children, we incorporate an open-source dataset comprising 10.43 hours of transcribed conversational speech on specific topics⁴¹. Additionally, we include VoIP, a dataset of telephone conversations (Alfano et al., 2014), within the same category.

B.10 Japanese

In addition to multilingual resources, our Japanese dataset includes educational content from Wikibooks⁴² and Wikijunior⁴³, as well as children’s books from Aozora Bunko⁴⁴.

Wikibooks is a collection of educational materials. From this dataset, we used the “Elementary School Learning” section, which targets Japanese elementary school students, typically aged 6 to 12. The content covers major school subjects, includ-

³⁷<https://children.logoslibrary.eu/>

³⁸<https://www.youtube.com/@CalimeroOfficial>

³⁹<https://www.invalsi.it/invalsi/index.php>

⁴⁰<https://www.mim.gov.it/>

⁴¹<https://magichub.com/datasets/>

⁴²<https://ja.wikibooks.org/wiki/>

⁴³<https://ja.wikibooks.org/wiki/Wikijunior>

⁴⁴<https://www.aozora.gr.jp/>

ing the Japanese language, social studies, mathematics, and science. We excluded pages that were still under construction, as well as those consisting primarily of numerical content (e.g., math drills). The resulting Wikibooks corpus contains approximately 0.2M words.

Wikijunior offers educational content designed for Japanese children aged approximately 8 to 11. As with Wikibooks, we excluded pages that were under construction or contained only numerical content. The final Wikijunior corpus consists of 75 pages, totaling approximately 0.07M words from Wikijunior.

Aozora Bunko is a Japanese digital library that provides access to literary works in the public domain. We used the aozorabunko-clean dataset⁴⁵, a cleaned version of the original collection, that includes only books whose copyrights have expired. This dataset contains storybooks, biographies, poetry, and other literary genres, with the majority being storybooks. It also contains Japanese translations of foreign literature. From this dataset, we selected only children’s books. The list of children’s book titles was scraped from the category-wise list of titles on Aozora Bunko⁴⁶. Books written in old character forms were excluded. This subset comprises 1,111 titles and totals approximately 8.7M words.

B.11 Persian

Our Persian dataset includes several curated subcategories designed to support both child-centered and educational language modeling. The final collection contains about 98.5 million words across 217,880 records and consists of four parts: Children’s Books, Educational Documents, Child-Directed Speech, and Subtitles used as supplementary padding.

Educational Documents. To construct this subset, we started with FineWeb2-HQ (Messmer et al., 2025), a high-quality, multilingual dataset built on top of FineWeb2 (Penedo et al., 2025), as the base for our educational subset. To identify educational content within the Persian subset, we fine-tuned an XLM-R (Conneau et al., 2020) model using a regression task inspired by the FineWeb-edu (Lozhkov et al., 2024) methodology. The training data for this model were annotated using Qwen2.5-

72B-Instruct (Qwen et al., 2025), following a 5-point additive rubric designed to assess the educational suitability of a document for primary to grade school learners. The documents were then scored between 0 and 5, which were later normalized to the 0–1 range. We trained the XLM-R model to predict these normalized scores. For our final dataset selection, we applied the trained model to Persian FineWeb2-HQ documents. We selected documents that (1) were under 3,000 words in length (to avoid structural drift across sections), (2) had a quality score of at least 0.35 based on FineWeb2-HQ metadata, and (3) received a predicted educational score of 0.9 or higher. This filtering ensures that the selected documents meet at least the first four points of the rubric, which means they are coherent, suitable for grade-school learners, and contain well-structured educational material.

Children’s Books. This subset includes child-friendly Persian texts sourced from two main corpora: Ririro (Persian section)⁴⁷ and GlotStoryBook.

Ririro. We scraped texts from the Persian section of the Ririro story collection. Although the content was mostly clean, we noticed minor inconsistencies in orthography and annotation. We applied light post-processing to fix punctuation, normalize spelling, and standardize diacritics, resulting in a clean and consistent corpus.

GlotStoryBook. This source included several very short entries, such as single-word or phrase-level records. To ensure data quality and narrative coherence, we filtered out all records with fewer than three words, resulting in the removal of 123 entries from an original total of 1,150. Notably, about one-third of the GlotStoryBook dataset consists of “fa-diacritics” texts, which are Persian sentences written with full diacritics. These fully vocalized texts are typically used in early literacy education in Persian-speaking contexts, particularly during the first stages of primary school. They are the first form of written Persian encountered by children as they begin learning to read and write, offering a bridge toward later reading of undiacritized Persian. Their inclusion enriches the dataset with pedagogically relevant material closely aligned with actual educational practice in early schooling.

⁴⁵<https://huggingface.co/datasets/globis-university/aozorabunko-clean>

⁴⁶<https://yozora.main.jp/>

⁴⁷<https://ririro.com/fa/>

Child-Directed Speech. We utilized Persian transcripts from the CHILDES project. Notably, although the spoken language is Persian, the transcripts are written in Latin script using phonetic representations, a transcription style known as Romanized Persian. Apart from standard normalization and deduplication, no further preprocessing was applied to preserve the phonetic and linguistic characteristics of child-directed speech.

Persian Subtitles. To meet our target budget of approximately 100 million words, we supplemented the dataset with Persian subtitle data. Subtitles were selected for their syntactic diversity and colloquial tone, helping to enrich the stylistic and lexical range of the dataset. The subtitles act as neutral padding rather than targeted educational or child-focused content.

B.12 Polish

In addition to the multilingual resources, we add three further data sources to the Polish data. The Wolne Lektury archive contains a large number of Polish ebooks. We systematically scraped all virtual bookshelves that contain child-directed/child-available literature and included all ebooks that could be plausibly encountered by children currently learning Polish. In order to do so, we consulted native speakers of Polish and articles on classical Polish children’s literature. Furthermore, we opt to include books that are translations of global children’s classics (e.g. *Tom Sawyer* or *Alice’s Adventures in Wonderland*), as they could plausibly encountered by children learning Polish. Besides these ebooks we also include all educational materials from the WikiJunior bookshelf of Polish Wikibooks, and educational materials from the Polish Wikikids website, which – despite its name – is not a classical wiki, but rather a general educational website.

Besides these child-directed resources, we were unable to find further child-available data. Unfortunately, no spoken Polish corpora are freely available. Although some larger Polish corpora exist, projects like the National Corpus of Polish only offer rudimentary search functions and no accessible data for our purposes.

B.13 Portuguese

Our BabyLM dataset for the Portuguese language consists primarily of sort stories from GlotStory-Books and Ririro, child-directed speech from the

Portuguese portion of CHILDES, and articles from a children’s wiki. We supplement this data with conversational spoken Brazilian Portuguese speech⁴⁸. Standard pre-processing steps were applied uniformly to the whole dataset. Finally, we note that during our initial collection efforts a variety of potentially relevant resources were found, but due to time constraints have not been included in this iteration of the data. Two such resources are PPORTAL, the Public Domain Portuguese-language Literature Dataset (Silva et al., 2021), and a collection of natural speech data from CORAA⁴⁹.

B.14 Romanian

The Romanian BabyLM corpus consists of texts from CHILDES (and more data added by Faiz ...). We additionally include data from two pre-existing resources. Chitez et al. (2024) introduce the LEMI Romanian children’s literature corpus, which consists of 33,154 words. We also include data the children’s portion of the Romanian Language Corpus collected by Midrigan Ciochina et al. (2020).⁵⁰ The corpus consists of children’s literature in its poetry and fairy tales section.

B.15 South African languages: Afrikaans, isiXhosa, isiZulu, Sesotho, Sepedi

The number of large-scale datasets and benchmarks for African languages has grown in recent years, but the African continent remains under-resourced and under-represented in NLP research (Ojo et al., 2025). Collecting BabyLM datasets for African languages presents several challenges. Besides lacking child-directed speech corpora, most African languages lack even domain-general datasets of sufficient quality and scale to approximate developmentally plausible training.

As a first attempt to create BabyLM datasets for African languages, we focus specifically on the linguistically diverse context of South Africa. South Africa has 12 official languages, some of which are commonly included in massively multilingual web-scraped datasets. Importantly, all languages have some high-quality, manually curated datasets that are publicly available. This is thanks to government initiatives, such as the South African Centre for Digital Language Resources

⁴⁸<https://magichub.com/datasets/brazilian-portuguese-conversational-speech-corpus/>

⁴⁹<https://sites.google.com/view/tarsila-c4ai/coraa-versions>

⁵⁰<https://lmidriganciochina.github.io/romaniancorpus/>

(SADiLaR)⁵¹, which prioritise the development of language resources in all official languages. After surveying available datasets across languages, we conclude that five languages are candidates for BabyLM datasets with meaningful amount of data: Afrikaans, isiXhosa, isiZulu, Sesotho (Southern Sotho), and Sepedi (Northern Sotho).

Afrikaans is comparably more resourced and we were able to collect a tier 2 corpus. For the other four languages, we were limited to tier 3 corpora. The proportion of data that is truly developmentally plausible varies between languages and, in some cases, falls short in comparison to higher-resourced languages. While limited in scale, our datasets demonstrate the practical feasibility of BabyLM research for low-resource languages. We hope our work serves as a starting point for future research on developmentally plausible language modelling for African languages.

Developmentally plausible data Only Afrikaans and Sesotho are represented in CHILDES, but we obtain child-directed and child-adult interaction corpora for all five languages from the SA-CDI project.⁵² We include children’s books for all five languages from the GlotStoryBook dataset (Kargaran et al., 2023), originally scraped from Nalibali⁵³, an initiative promoting children’s literacy in South Africa. For educational content, we include high school exams (Sibeko and Zaanen, 2023) for language subjects (home language and first additional language) for all five languages and QED (Abdelali et al., 2014) for Afrikaans, isiXhosa, and isiZulu. For isiXhosa, we include the descriptive sentences in T2X (Meyer and Buys, 2024), a data-to-text dataset containing simplified isiXhosa sentences describing (subject, relation, object) triples in a knowledge base.

Padding corpora To match the target dataset sizes (tier 2 Afrikaans and tier 3 for isiXhosa, isiZulu, Sesotho, and Sepedi), we include additional high-quality data to supplement the developmentally plausible data as needed. For Afrikaans, we include OpenSubtitles (Lison and Tiedemann, 2016). For all five languages we include language-specific Wikipedia corpora. This still leaves us short for Sesotho and Sepedi. For Sepedi, we include government news articles from Vukenzele (Lastrucci et al., 2023). Finally, we use sentences

from parallel corpora for machine translation to reach our target sizes for Sesotho and Sepedi. For Sepedi we include the highest quality Sepedi sentences in WMT22 (Adelani et al., 2022), as measured by language identification score. For Sesotho we include sentences from the Autshumato English-Sesotho Parallel Corpus (McKellar, 2022).

B.16 Indonesian and its local languages: Javanese, Sundanese, Balinese, Buginese, Makassarese, Minangkabau, Acehnese

Recent years have seen a significant increase in resources for Indonesian and its local languages, mainly due to collective efforts by NusaCrowd (Cahyawijaya et al., 2023a) and SEACrowd (Love-nia et al., 2024). These initiatives have contributed a wide range of datasets, including conversational corpora, written texts, and multilingual collections. However, developmentally plausible and child-related data are still lacking. We can only find one dataset available from those collective efforts: ASR-INDOCSC, which consists of 4.5 hours of daily conversational speech from children in Indonesia, along with multilingual resources.

The main sources for cognitively and developmentally plausible data for Indonesian and its local languages come mainly from books obtained from a repository provided by the Ministry of Education & Culture⁵⁴. These are primarily educational books and storybooks for children aged 2 to 12. Since these books are in PDF format, we used PyPDF2 (Fenniak et al., 2022) and Tesseract (Smith, 2007) to extract their content. For data preprocessing, we use Gemma3-27B (Team et al., 2025) for content filtering in three steps: filter out non-child-related books, clean and reformat the extracted book content, and then remove non-child-related content. After cleaning, GlotLID v3 (Kargaran et al., 2023) was used for language detection and grouping, allowing data collection for Javanese, Sundanese, Balinese, Buginese, Makassarese, Minangkabau, and Acehnese. Another major source is the Bobo children’s magazine⁵⁵, which contains child-targeted articles from January 2020 to May 2025, all of which are exclusively in Indonesian. In addition to these, we incorporated data from multilingual resources, specifically GlotStoryBook (Kargaran et al., 2023) and Ririro⁵⁶, for Indonesian language data.

⁵¹<https://repo.sadilar.org/>

⁵²<https://sa-cdi.org/>

⁵³<https://nalibali.org/>

⁵⁴<https://repositori.kemdikbud.go.id>

⁵⁵<https://bobo.grid.id>

⁵⁶<https://ririro.com>

To pad the data and reach the required tiers, OpenSubtitles (Lison and Tiedemann, 2016) data were utilized for Indonesian to reach Tier 1. For local languages to reach Tier 3, we prioritized high-quality, manually curated datasets from NusaX (Winata et al., 2023), NusaWrites (Cahyawijaya et al., 2023b), and NusaDialogue (Purwarianti et al., 2025), followed by Wikipedia and MADLAD-400 (Kudugunta et al., 2023) data for additional padding as needed.

B.17 Spanish

As the predominant language in 21 countries, Spanish is a pluricentric language and exhibits rich diatopic variations. Far from being a homogeneous language, it encompasses a wide range of national and regional varieties, marked by distinct morphosyntactic and lexical features (Mayor-Rocher et al., 2025). As such, the term “Spanish” does not denote a single standardized form, but rather a set of linguistic norms shaped by diverse cultural and geographic contexts. The resources compiled in this dataset reflect this inherent diversity: our search for developmentally plausible materials was deliberately international, resulting in the inclusion of content from at least eight different countries.

Children’s Books A substantial portion of children’s books is sourced from the Elejandría collection⁵⁷, which features 19 translated bedtime stories from classical authors like Andersen, Grimm, and Perrault; 20 translated young adult classics, including “Gulliver’s Travels” and “Alice in Wonderland”; and 35 original Spanish-language books by authors from Spain, Uruguay, Mexico, Nicaragua, Cuba, and Argentina, categorized under Discovering Spain and Hispanic American Literature. Additional books were sourced from the Logos Group library, which granted us access upon request. This collection includes Spanish translations of well-known children’s literature, such as The Adventures of Tom Sawyer, as well as a smaller number of original Spanish texts. It also features songs, traditional Christmas carols, legends, and famous fables like The Ant and the Grasshopper. Our dataset also includes a range of children’s stories, fairy tales, poems, traditional literature, and songs accessed via the Ministries of Education of Argentina

⁵⁷<https://www.elejandria.com/coleccion/>

⁵⁸ and Colombia⁵⁹, the provincial government of Salta in Argentina⁶⁰, and the educational website *educ.ar*⁶¹.

Child-Accessible Speech To capture spoken Spanish that is accessible to children, we incorporated two complementary resources. First, we included an open-source dataset from MagicHub⁶², comprising 5.56 hours of transcribed conversational speech in Peninsular Spanish. This dataset features 17 dialogues recorded between four pairs of speakers, covering a variety of everyday topics. Additionally, we incorporated the SpinTX video archive⁶³, which offers curated video clips and transcripts from the Spanish in Texas Corpus. This collection of interviews with bilingual Spanish speakers residing in Texas covers a wide range of topics relevant to daily life, including family, friendship, food, culture, parenting, education, and school.

B.18 Ukrainian

The Ukrainian dataset is a collection of different resources. To the best of our knowledge, there is no CHILDES-like corpus for the Ukrainian language; therefore, it has been substituted with a set of monolingual and multilingual data.

Developmentally plausible data. For the majority of developmentally plausible data, we use the GRAC corpus (Shvedova and Lukashevskyi, 2024). This corpus consists of copyright-free texts concerning Ukraine till 1954. The dataset is heavily filtered, reducing from 100M tokens to 29M, to extract the most developmentally plausible data. First, language filtering restricts content to Ukrainian, excluding all other languages, including English, German, Russian, and others. Style-based filtering removes journalistic content, personal memoirs, religious materials, public speeches, official documents, and texts with unknown style classifications. Additionally, non-fiction works published before 1900 are excluded to maintain temporal relevance. The remaining texts are categorized into educational content (academic materials and popular science works), child-appropriate books (fiction,

⁵⁸<https://www.argentina.gob.ar/educacion/historiasxleer>

⁵⁹<https://v1.maguared.gov.co/serie-leer-es-mi-cuento-todos-los-titulos/>

⁶⁰<https://planeamiento.edusalta.gov.ar/>

⁶¹<https://www.educ.ar/>

⁶²<https://magichub.com/datasets/>

⁶³<https://spintx.org/>

folklore, and poetry), and other materials (internet communication and private oral content). Additionally, we utilize the Ukrainian portion of Wikisource (Wikimedia Foundation, 2025) as a source of fairy tales and fiction books, thereby expanding the dataset by an additional 1 million tokens.

Padding corpora. To expand the developmentally plausible data, we incorporate the previously mentioned GlotStorybook and Ririro datasets. Wikipedia serves as a significant source of encyclopedic content, contributing approximately 29.1M tokens. The FineWeb-C corpus provides an additional 174K tokens of contemporary language use. Finally, OpenSubtitles contributes nearly 29.5M tokens of conversational Ukrainian text from movie and television subtitles, to which a child would most likely be exposed.

B.19 Other Languages

For the rest of the languages in the BabyBabelLM dataset, no language-specific resources were collected. Instead, these languages are populated by multilingual data resources, namely: CHILDES, GlotStoryBooks, Ririro, and Child Wikis. These languages are: *Basque, Croatian, Czech, Danish, Estonian, Hebrew, Hungarian, Icelandic, Korean, Norwegian, Romanian, Russian, Serbian, Turkish, Swedish, and Welsh* for a total of 16 out of 45 languages. We welcome contributions for these, and other languages, details presented in our GitHub.

Language	ISO 639-3	Tier	Byte Premium	Actual Data Size (MB)	Transcription Tokens	Education Tokens	Books, Wiki, News Tokens	Subtitles Tokens	Padding Tokens	Total Tokens
Chinese	zho	Tier 1	0.989	537.15	15,202,509	13,305,706	16,493,430	0	105,307,983	150,309,628
French	fra	Tier 1	1.174	634.88	6,234,743	0	13,987,611	106,358,431	0	126,580,785
Bulgarian	bul	Tier 1	1.812	981.07	143,293	0	24,799,312	90,435,735	0	115,378,340
Indonesian	ind	Tier 1	1.179	638.87	17,824	62,188	17,225,662	96,044,750	0	113,350,424
Dutch	nld	Tier 1	1.052	569.59	3,304,756	19,146,045	17,428,015	70,006,748	0	109,885,564
German	deu	Tier 1	1.054	569.1	8,518,785	257,233	7,655,975	91,478,846	0	107,910,839
English	eng	Tier 1	1.000	539.59	36,814,704	0	41,357,314	19,699,367	1,068,358	98,939,743
Persian	fas	Tier 1	1.597	867.4	234,221	94,320,938	67,165	3,915,679	0	98,538,003
Ukrainian	ukr	Tier 1	1.751	945.57	0	12,003,085	16,786,100	29,496,904	29,307,158	87,593,247
Serbian	srp	Tier 2	1.425	77.25	1,489,908	0	29,896	13,707,246	0	15,227,050
Cantonese	yue	Tier 2	0.862	46.74	2,982,684	0	191,861	0	11,870,650	15,045,195
Japanese	jpn	Tier 2	1.322	71.79	3,656,514	291,053	9,712,521	721,565	0	14,381,653
Portuguese	por	Tier 2	1.098	59.5	956,441	0	382,562	10,348,599	0	11,687,602
Swedish	swe	Tier 2	1.021	55.21	750,286	0	526,330	9,810,043	0	11,086,659
Greek	ell	Tier 2	1.967	106.74	2,123,853	2,577,703	1,402,782	4,847,444	0	10,951,782
Polish	pol	Tier 2	1.077	58.29	1,257,155	0	48,831	8,906,729	0	10,212,715
Estonian	est	Tier 2	0.968	52.37	1,026,491	0	0	8,814,184	0	9,840,675
Spanish	spa	Tier 2	1.084	58.75	2,978,384	0	5,385,855	1,344,853	0	9,709,092
Italian	ita	Tier 2	1.067	57.96	1,189,631	990,522	6,797,154	471,404	18,976	9,467,687
Afrikaans	afr	Tier 2	1.037	56.28	272,442	116,380	153,914	1,315,741	7,448,952	9,307,429
Welsh	cym	Tier 2	1.027	55.39	1,109,683	0	0	405,811	7,196,648	8,712,142
Arabic	ara	Tier 2	1.465	79.61	3,304,730	0	1,672,594	3,349,526	0	8,326,850
Basque	eus	Tier 2	1.06	57.06	201,402	0	1,716,026	3,176,681	3,095,188	8,189,297
Hebrew	heb	Tier 2	1.356	73.38	1,045,939	0	0	0	6,117,225	7,163,164
Sesotho	sot	Tier 3	1.166	6.31	420,926	106,253	141,012	0	557,317	1,225,508
Sepedi	nso	Tier 3	1.116	6.06	0	92,589	122,083	0	871,565	1,086,237
Buginese	bug	Tier 3	1.228	6.67	0	0	41,174	0	961,405	1,002,579
Romanian	ron	Tier 3	1.115	6.1	294,696	0	284,101	393,308	0	972,105
Acehnese	ace	Tier 3	1.242	6.74	0	0	242,613	0	725,581	968,194
Javanese	jav	Tier 3	1.147	6.23	0	0	307,282	0	645,365	952,647
Balinese	ban	Tier 3	1.27	6.87	0	0	63,826	0	874,899	938,725
Icelandic	isl	Tier 3	1.154	6.27	452,099	0	0	470,031	0	922,130
Croatian	hrv	Tier 3	0.99	5.39	469,078	0	0	445,976	0	915,054
Makassarese	mak	Tier 3	1.251	6.79	0	0	34,080	0	873,230	907,310
Norwegian	nor	Tier 3	1.125	6.11	404,670	0	290	0	496,473	901,433
Sundanese	sun	Tier 3	1.097	5.96	0	177	17,264	0	874,647	892,088
Danish	dan	Tier 3	1.021	5.53	372,836	0	8,848	0	457,152	838,836
Russian	rus	Tier 3	1.823	9.88	151,640	0	92,462	0	561,542	805,644
Korean	kor	Tier 3	1.293	7.07	649,349	0	6,794	138,940	0	795,083
Minangkabau	min	Tier 3	0.95	5.16	0	0	122,536	0	663,669	786,205
Czech	ces	Tier 3	1.036	5.64	377,313	0	0	0	385,263	762,576
isiZulu	zul	Tier 3	1.164	6.31	62,772	56,641	96,383	5,402	532,402	753,600
Hungarian	hun	Tier 3	1.02	5.55	391,041	0	6,234	0	322,636	719,911
isiXhosa	xho	Tier 3	1.199	6.52	74,950	65,208	98,144	29,099	412,004	679,405
Turkish	tur	Tier 3	1.044	5.72	248,397	0	11,193	0	405,478	665,068

Table 2: Detailed data statistics for all languages in the BabyBabelLM dataset. Tiers indicate target size equivalence to English tokens: Tier 1 (100M), Tier 2 (10M), Tier 3 (1M).

Field	Type	Values	Description
text	string	<i>Una volta, c'erano...</i>	The content of the document.
category	string	Transcription <ul style="list-style-type: none"> child-directed-speech child-available-speech Education <ul style="list-style-type: none"> educational Wiki, News, Books <ul style="list-style-type: none"> child-books child-wiki child-news Subtitles <ul style="list-style-type: none"> subtitles qed Padding <ul style="list-style-type: none"> padding-wikipedia padding-[placeholder] 	Speech directed to children and speech produced by children. Speech children are exposed to without being the main recipients (e.g., adult conversations). School textbooks, exams, and other educational material designed for children. Books and stories created for children. Children wiki articles. News directed to children. Subtitles for child-appropriate material (e.g., children TV shows). Subtitles from the QED dataset. Wikipedia articles. Other forms of padding, used primarily for low-resource languages.
data-source	string	CHILDES, www.ririro.com, ...	The source of the document: dataset name, url, or itwm identifier.
script	string	Latn, Grek, Arab, ...	The script of the text: a validated ISO-15924 code string.
age-estimate	string	3-6, children, adults, n/a, ...	For text data: estimated age of target audience. For speech data: estimated age of speakers.
license	string	cc-by-nc, public domain, ...	The license under which the document is released.
misc	string	{"info": "...", ...}	Optionally included supplementary information, as a valid JSON string.
num-tokens	integer	15364	The number of white-separated (or tokenizer-based) tokens present in the document text.
doc-id	string	7a2b3a1d9...	Unique document ID computed as a sha256 string of it's content, used for de-duplication.

Table 3: Document-level schema for the BabyBabelLM datasets. For each document field we define it's type, include sample values, and give a description of it's use and contents.

	Language	multihimp	monohimp	include	bnaiama	multitli	belebele	global-mm1u	arc	hellaswag	xnli	xcopa	xstorycloze	xcomps	wino grande	truthfulqa	sh200
TIER 1	Random	50.0	50.0	25.0	25.0	33.3	33.3	25.0	25.0	25.0	33.3	50.0	50.0	50.0	50.0	25.0	25.0
	Standard Arabic	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
	Bulgarian	83.5	—	37.5	33.7	57.4	50.6	—	32.8	29.7	37.1	—	52.8	—	51.0	30.8	25.0
	Chinese	—	83.0	52.5	53.3	54.6	61.7	45.8	51.5	37.1	35.5	63.8	63.3	64.5	56.0	26.0	25.0
	Dutch	82.8	—	33.4	40.4	43.0	50.1	37.6	33.2	30.0	—	—	54.0	53.2	51.6	24.0	25.0
	English	98.2	81.0	—	66.5	58.1	55.9	50.0	57.3	40.0	53.4	73.0	67.0	—	55.6	22.4	27.9
	French	96.7	84.6	42.2	46.7	61.0	56.7	38.9	41.4	33.4	44.7	—	58.1	53.0	54.2	26.0	25.5
	German	93.7	86.8	36.0	51.9	33.2	54.0	38.7	37.4	32.3	44.0	—	56.3	54.9	53.4	26.9	25.5
	Indonesian	—	—	44.9	46.3	36.0	53.1	38.8	38.8	32.3	—	57.8	55.7	—	55.8	23.3	25.0
	Persian	81.1	—	33.6	27.9	49.9	46.0	32.6	31.4	29.3	—	—	54.8	52.0	52.8	24.3	25.0
TIER 2	Ukrainian	85.1	—	47.5	33.4	32.7	45.0	34.7	34.4	30.3	—	—	50.9	52.8	52.5	27.6	25.0
	Afrikaans	—	—	—	34.1	37.7	43.8	—	28.3	28.3	—	—	51.2	—	51.4	21.8	25.0
	Basque	89.0	41.6	30.0	—	—	33.9	—	25.4	—	33.5	49.6	49.6	—	—	—	25.0
	Estonian	61.5	—	30.8	23.8	38.1	35.7	—	26.5	27.4	—	47.4	49.3	—	49.3	22.6	25.0
	Greek	92.5	—	33.1	30.7	35.9	42.7	32.5	29.2	29.2	36.3	—	51.9	50.6	52.4	28.4	25.0
	Hebrew	73.0	64.4	42.7	28.3	45.5	40.3	30.9	30.6	28.6	—	—	51.5	50.5	49.2	31.5	25.5
	Italian	88.4	—	45.8	46.4	38.6	53.6	40.2	39.5	33.2	—	56.8	57.0	—	52.4	29.4	25.5
	Japanese	—	74.6	44.7	34.1	58.0	47.1	38.1	39.1	31.4	—	—	56.2	52.8	51.2	27.2	26.0
	Polish	81.5	—	38.1	31.0	39.7	48.6	36.3	34.7	29.8	—	—	53.8	—	53.5	27.6	25.0
	Portuguese	93.4	—	43.2	41.5	38.6	55.2	36.7	41.2	33.9	—	—	58.7	—	53.1	25.5	25.5
	Serbian	—	—	29.3	31.4	41.5	47.8	32.4	28.5	29.6	—	—	53.5	—	52.3	29.2	25.0
	Spanish	93.2	—	42.2	47.7	61.8	48.7	37.6	42.2	34.5	41.4	59.6	58.0	54.3	53.9	28.7	25.5
	Swedish	—	—	—	42.9	53.6	49.3	36.0	32.1	30.3	—	—	52.2	—	50.3	25.7	25.5
	Welsh	70.8	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
	Yue Chinese	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	26.0
TIER 3	Achinese	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	25.0
	Balinese	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	25.0
	Buginese	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	25.0
	Croatian	—	—	33.3	32.1	35.9	46.9	—	30.6	29.2	—	—	51.6	—	52.3	28.2	25.0
	Czech	77.8	—	—	30.1	38.8	48.9	36.0	33.5	29.7	—	—	53.1	—	50.8	24.7	25.0
	Danish	92.0	—	—	40.2	36.8	46.4	—	32.0	30.0	—	—	53.8	—	52.4	27.9	25.5
	Hungarian	84.5	—	32.5	26.6	53.1	40.7	—	29.3	28.4	—	—	52.5	51.5	51.1	27.4	25.0
	Icelandic	63.2	—	—	21.7	35.8	35.4	—	26.0	26.7	—	—	46.5	—	49.0	21.9	25.0
	Javanese	—	—	—	36.8	37.0	36.3	—	28.4	27.9	—	—	49.9	—	50.6	19.1	25.0
	Korean	—	—	40.6	32.3	53.1	47.4	35.1	38.0	30.2	—	—	54.3	53.2	51.1	26.0	30.4
	Makasar	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
	Minangkabau	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	25.0
	Norwegian	—	—	—	40.7	38.5	46.3	—	31.8	29.7	—	—	53.2	—	50.9	24.0	25.0
	Pedi	—	—	—	—	—	29.6	—	26.1	—	—	—	—	—	—	—	25.0
	Romanian	85.0	—	—	35.4	56.0	51.7	37.3	33.2	30.5	—	—	53.4	—	53.2	25.5	25.0
	Russian	89.6	85.6	44.4	35.9	49.2	53.4	37.9	41.8	33.5	43.7	—	59.4	55.9	53.7	27.4	25.0
	Southern Sotho	—	—	—	—	—	29.1	27.6	—	—	33.3	—	—	—	—	—	25.0
	Sundanese	—	—	—	—	—	31.6	—	—	—	—	—	—	—	—	—	25.0
	Turkish	79.2	75.4	38.5	32.7	49.0	41.8	34.1	36.3	29.4	38.5	56.2	52.3	51.3	48.6	26.9	25.0
	Xhosa	—	—	—	—	—	28.7	26.2	—	—	33.3	—	—	—	—	—	25.0
	Zulu	—	—	—	—	—	32.0	25.0	26.3	—	33.3	—	—	—	—	—	25.0

Table 4: Performance of the Qwen3-0.6B trained on BabyBabelLM. All scores denote average 0-shot accuracy. Columns are sorted by the column order of Table 1.