CurLL: Curriculum Learning of Language Models

Anonymous Author(s)

Affiliation Address email

Abstract

We introduce a comprehensive continual learning dataset and benchmark (CurlL) grounded in human developmental trajectories from ages 5-10, enabling fine-2 grained assessment of models' ability to acquire new skills. CURLL spans five 3 developmental stages (0-4) covering ages 5-10, supported by a skill graph that 4 breaks down broad skills into smaller abilities, concrete goals, and measurable 5 indicators, while also capturing which abilities build on others. We generate a 23.4B-token synthetic dataset with controlled skill progression, vocabulary complexity, and format diversity, comprising paragraphs, comprehension-based QA 8 (CQA), skill-testing QA (CSQA), and instruction-response (IR) pairs. Stage-wise 9 token counts range from 2.12B to 6.78B tokens, supporting precise analysis of 10 forgetting, forward and backward transfer. Using a 135M-parameter transformer 11 trained under independent, joint, and sequential setups, we show trade-offs in skill 12 retention and transfer efficiency. 13

1 Introduction

The capacity for lifelong learning in humans is not just a practical advantage but a fundamental 15 aspect of intelligence itself [Kudithipudi et al., 2022, Yan et al., 2024, Schmidgall et al., 2023]. The 16 continual learning (CL) problem thus is one of the grand challenges for achieving human-like artificial 17 intelligence. It addresses the core problem of how computational systems can progressively acquire, integrate, and refine knowledge over extended periods without compromising earlier capabilities. 19 For language models (LMs), this challenge is particularly interesting: despite their impressive 20 performance across various tasks, these models face a fundamental limitation in that their skill-set 21 and knowledge of the world become static after training, frozen at the point of deployment [Shi 22 et al., 2024, Wu et al., 2024, Bell et al., 2025]. Despite the importance of the CL problem for LMs, 23 current evaluation methodologies suffer from significant limitations: 1) Poor skill control: Existing 24 25 benchmarks often lack precise control over the specific skills being tested, making it difficult to isolate the effects of learning new capabilities [Liu et al., 2025, Rivera et al., 2022]. 2) Unclear knowledge dependencies: The relationships between skills are rarely explicitly modeled, missing 27 out on important transfer effects [Zheng et al., 2025, Nekoei et al., 2021]. 3) Inadequate forgetting 28 metrics: Many evaluations fail to properly measure catastrophic forgetting across sequential learning 29 tasks [Chen et al., 2023a, Huang et al., 2023]. 30

To address these gaps, we introduce a dataset (Curll) to train and evaluate continual learning algorithms for language models. Coming up with a set of skills with a rich structure and dependencies is a challenge in the construction of such a dataset. We find such a source of skills in human education. Curll is grounded in the curriculum for human education from ages 5–10, divided into five developmental stages (0–4). Each of these stages represent one human-year. Our framework incorporates 1,300+ fine-grained skills with dependencies codified in a skill graph having skills as nodes with the edges capturing a prerequisite relationship. The edges are weighted on a scale of (1–5) to capture dependency strength. Starting from this set of skills, we generate a synthetic

dataset of 23.4B tokens, with controlled vocabulary complexity (stage-specific word sampling from Age-of-Acquisition data as seed) and multiple formats (paragraphs, comprehension QA, skill-testing QA, instruction-response). Each stage's dataset ranges from 2.12B to 6.78B tokens, enabling fine-grained evaluation at indicator, skill, and stage levels. Our code, dataset (stages 0–4), and skill graph will be publicly released. Our contributions include: a) The idea of grounding skills in human education curriculum in the context of CL. b) A synthetic data generation pipeline spanning 5 developmental stages with stage-specific vocabulary and explicit skill dependencies. This pipeline gives us a benchmark with fine-grained control over measuring skill transfer, forgetting and sample efficiency c) A skill graph-based dependency model that explicitly captures prerequisite relationships between learning objectives, enabling nuanced analysis of skill transfer and forgetting.

9 2 Related Work

Many datasets and benchmarks exist for continual learning of LMs [Jang et al., 2021, Li et al., 2025]. 50 TRACE [Wang et al., 2023] highlights that existing benchmarks are too simple or are already included 51 in instruction-tuning sets. MMLM-CL [Zhao et al., 2025] notes the limited real world applicability in 52 benchmarks. OCKL [Wu et al., 2023] proposes new metrics for measuring knowledge acquisition rate and knowledge gap but concentrates on knowledge-intensive tasks as compared to procedural tasks. 54 TemporalWiki [Jang et al., 2022] is for updating factual information in LMs based on temporal data. 55 SuperNI contains a variety of traditional NLP tasks and serves as a practical benchmark for continual 56 learning of large language models [He et al., 2024]. Despite these developments, these benchmarks 57 are often considered unsuitable for evaluating state-of-the-art LMs [Wang et al., 2023, Razdaibiedina 58 et al., 2023, Scialom et al., 2022, Zhang et al., 2015]. These benchmarks often emphasize artificial 59 task boundaries He et al. [2024], lack temporal and distributional complexity. Moreover, these 60 datasets do not offer precise control over skills or information to validate the effectiveness of existing 61 solutions for continual learning. Skill-it [Chen et al., 2023b] emphasizes the problem but only 62 introduces a data sampling algorithm for continual pretraining by arranging the skills in a increasing 63 order of complexity. Other existing works [Khetarpal et al., 2020, Greco et al., 2019, Xu et al., 64 2024] discuss the importance of skill distinction and its effect on evaluating continual learning. In 65 contrast, our work is grounded in human developmental curricula and enables fine-grained evaluation 66 of transfer, forgetting, and sample efficiency beyond what existing benchmarks support.

68 3 Dataset Setup

77

78

79

80

81

84

85

86

87

Our framework is grounded in human learning curriculum, with the dataset designed to mimic the 69 developmental stages from age 5-10. We use two established educational frameworks to develop 70 our skill taxonomy: the Early learning Outcomes framework (ELOF) for children aged 5¹ and the 71 Cambridge curriculum for children aged 5-10². These frameworks help us define fine grained notion 72 of skills as specified by a skill-tuple that consists of four components: 1) Skills³: High-level domains 73 or subjects (e.g. Mathematics, Science). 2) Sub-skills: Specific components within a skill (e.g., 74 Counting and Cardinality). 3) Goals: Broad statement of learning expectations within a sub-skill. 4) 75 Indicators: Specific, observable behaviors that demonstrate mastery of a goal. 76

The ELOF framework has five broad areas: Approaches to Learning, Social and Emotional Development, Language and Literacy, Cognition, and Perceptual, Motor, and Physical Development. Cambridge Primary Curriculum covers subjects including English, Mathematics, Science, Computing, and Global Perspectives. The curriculum structure flows from subjects (renamed as skills in our framework) to domains/strands (renamed as subskills), then to substrands (goals), each with specific learning objectives (indicators). We also adopt the notion of stages from the Cambridge curriculum in our framework, where each stage corresponds to one year starting from age 5. Therefore, we have 5 stages in our framework, where stage 0 denotes ages up to 5, stage 1 denotes age 5-6 and so on. The number of skill-tuples in our framework is the same as the number of indicators present in stages 0-4, statistics of which are mentioned in Table 1. We construct a skill graph, which is a directed graph that has indicators as nodes, with edges representing prerequisite relationships weighted from 1-5 to indicate dependency strength. These edges model how skills are built on each other in developmental

¹U.S. Department of Health & Human Services, Administration for Children & Families [2024]

²Cambridge Assessment International Education [2025]

³"Skill" here has a specific meaning, which is different from the general notion of skill used before

| Table 1: Dataset statistics across developmental stages (0–4), including | ding total tokens |
|--|-------------------|
|--|-------------------|

| Stage | | Skills | & Goals | | | Instances | | |
|-------|----------|--------------|---------|--------------|-------|-----------|------------|-------|
| | # Skills | # Sub-skills | # Goals | # Indicators | # CQA | # CSQA | # IR Pairs | in Bn |
| 0 | 7 | 24 | 59 | 182 | 1.0M | 3.01M | 3.30M | 2.12 |
| 1 | 7 | 29 | 86 | 292 | 20.2M | 4.04M | 4.10M | 3.47 |
| 2 | 6 | 26 | 67 | 249 | 23.5M | 4.70M | 4.78M | 4.56 |
| 3 | 6 | 26 | 68 | 271 | 31.2M | 6.24M | 6.29M | 6.47 |
| 4 | 6 | 23 | 70 | 349 | 27.4M | 5.49M | 5.52M | 6.78 |

stages⁴. While the skill graph isn't directly used for skill data generation, it provides insights for analyzing continual learning patterns and interpreting evaluation results (see Appendix A).

3.1 Synthetic Data Generation

qq

Our synthetic data consists of instances, each mimicking a situation a child might encounter. Instances are of three types: (1) IR: an instruction-response pair, where the instruction is about some general world knowledge, (2) CQA: context-based question-answers for testing comprehension, (3) CSQA: context-based question-answers for testing skills. A context is a short piece of text which forms the basis of the corresponding question-answer pairs in the instance. Contexts can be of multiple types as specified by a template: e.g., a simple narrative, or a dialogue. IR-pairs can also have different types specified by templates, e.g., mimic action or follow simple direction (Appendix A for examples).

Instances are generated by prompting an LLM with a *seed*. A seed consists of a skill-tuple, vocabulary seed, instance type, template. This choice is crucial for ensuring diversity and coverage of our data. This tuple is also our way to ground the generations in the skill graph. To generate one instance of the data, we first construct a seed: each skill-tuple is combined with a vocabulary seed for that stage, an instance type and a template for that instance type. If the instance type is CQA or CSQA, then we first generate the context, and then using the context, we generate the corresponding question-answers. If the instance type is IR then we directly generate the instruction-response pairs. The prompts for all the generations and details of the generation process are presented in Appendix A. In our dataset, each instance includes the seed used to generate it as part of its metadata. We generated data for stages 0-4, containing a total of 23.4B tokens (Table 1).

We measure diversity of generated data using: 1) Diversity as reciprocal of compression ratio using gzip Gailly and Adler [1992]. 2) The intra- and inter-text deduplication rate as calculated by semantic deduplication. Cross-stage analysis shows higher diversity and lower deduplication rate (<5%) between stages compared to intra-stage results, confirming that content evolves meaningfully across developmental progression while maintaining stage-specific uniqueness. See Appendix B for more details. We also measure progression in the difficulty of the skills as the stage number increases. We sample 500K instances from each stage for each data type and run statistical readability tests. Means across multiple readability metrics are reported in Appendix B.3. The readability tests show that as stages progress, the texts also become increasingly challenging. At least 50 random instances from each dataset per stage were manually analysed, revealing that CQA data for all stages was found to be accurate. IR and CSQA data had certain patterns like excessive use of discourse markers for early stages and verbose response to instructions. We choose 25 instances per indicator for test set resulting in 5k-7k samples per stage. Since the data is synthetically generated at scale, we reserve the highest quality samples for the test set. 100 instances per indicator instance type are sampled randomly and rated by LLM on a scale 1-5. Top 25 instances are selected for test, next 25 for validation set.

4 Experiments and Results

We conduct preliminary experiments to validate that the dataset exposes meaningful challenges: whether models can retain earlier-learned skills, how sequential training affects generalization, and to what extent transfer across related skills occurs. By analyzing these at the granularity of skills, we demonstrate that Curll enables insights that are not visible in existing benchmarks. Unlike

⁴an LLM (Gemma3-27B-IT is used for all LLM inferences throughout this work) is used to predict the edges ⁵Uses pre-defined words to predict the grade of a text (https://github.com/cdimascio/py-readability-metrics)

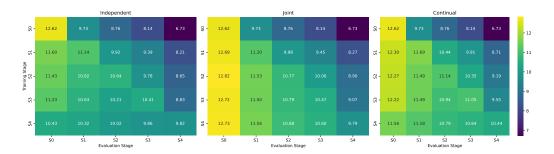


Figure 1: Stage-wise results for all training setups. Independent corresponds to models trained on a single stage, joint to models trained on mixtures of data up to a stage, and continual to sequential upto a stage. Heatmaps report summed correctness scores across all test formats (IR, CQA, CSQA)

traditional language model training that includes two stages: pretraining and then finetuning, we do a single phase training. All instance types i.e. CQA, CSQA, and IR are included in the same phase. Since all of them are question-answers, with and without context paragraphs, we use a standard chat template to train the language models from scratch. Smollm2-135M parameter model is used as the base architecture. All training runs are performed on one full epoch of the data. Learning rate of 5e-3 and effective batch size of 1536 instances remain unchanged across experiments. We use a context length of 1024. Other training and inference related hyper-parameters are mentioned in the Appendix C. We perform three types of training: 1) Independent (M_i) : The model is trained from scratch on data of each stage independently 2) Joint (M_{ij}) : Jointly trained on a mixture of stages. The data from different stages is combined and shuffled randomly. 3) Continual (M_{i-j}) : The model is first trained on stage i, then stage j, then stage k and so on. To evaluate the trained models, the instances from test set are passed through the chat template and the model is asked to complete the generation post instruction. These inferences along with the prompt is passed to an LLM to rate on a scale of 1-5. This is followed for all three types of test sets. Each model is evaluated on test sets of all stages. The main objective of the rating is to evaluate the correctness of the model inference with some weightage to the stage on which the model is being evaluated. The summation of scores across test set types (IR, CQA, CSQA) is presented in Figure 1. The individual scores, prompts and rubrics for evaluation are available in the Appendix E.

Joint models (M_{ij}) generalize better to later stages and maintain strong performance on trained stages compared to independent models (M_i) . Continual models (M_{i-j}) , however, achieve the best performance on later stages but suffer degradation on earlier ones. Sequential (continual) ordering improves generalization but also induces forgetting of earlier skills, which is counter-intuitive since later skills depend on foundational ones. The skill graph helps explain this. The largest performance gaps between joint and continual training occur for "Perceptual, Motor, and Physical Development" and "Digital Literacy". Both have very few outgoing edges in the skill graph (Appendix D), meaning their indicators rarely serve as prerequisites for later skills.

5 Conclusion

We introduced (CurlL), a novel continual learning evaluation framework for language models grounded in human developmental curricula. (CurlL) combines a directed, weighted skill graph of over 1,300 fine-grained skills with a 23.4B-token synthetic dataset that controls stage-wise vocabulary, difficulty, and format. The skill graph serves as a diagnostic tool: its metadata enables fine-grained control over the number of instances and skills seen during training, supports evaluation of sample efficiency, and allows targeted testing of transfer effects (e.g., whether learning Skill A improves Skill B). Forgetting, forward transfer, backward transfer, and data efficiency can all be measured at the levels of skills, sub-skills, and indicators. This enables richer analysis than stage- or task-level metrics in existing benchmarks, which typically report only overall accuracy on entire tasks (e.g., classification or QA) without revealing which underlying abilities are gained or lost. Our experiments with independent, joint, and sequential training demonstrate that simply changing the order of data presentation affects both generalization and forgetting. Finally, the scalable data generation pipeline enables exploring continual pretraining in a controlled yet realistic setting.

69 References

- 170 Jack Bell, Luigi Quarantiello, Eric Nuertey Coleman, Lanpei Li, Malio Li, Mauro Madeddu, Elia
- Piccoli, and Vincenzo Lomonaco. The future of continual learning in the era of foundation models:
- Three key directions. ArXiv, abs/2506.03320, 2025. URL https://api.semanticscholar.
- org/CorpusId:279155222.
- 174 Cambridge Assessment International Education. International curricu-
- lum. https://www.cambridgeinternational.org/why-choose-us/
- benefits-of-a-cambridge-education/international-curriculum/,
- 177 2025. URL https://www.cambridgeinternational.org/why-choose-us/
- benefits-of-a-cambridge-education/international-curriculum/. Accessed
- 179 2025-08-20.
- Ernie Chang, Matteo Paltenghi, Yang Li, Pin-Jie Lin, Changsheng Zhao, Patrick Huber, Zechun Liu,
- Rastislav Rabatin, Yangyang Shi, and Vikas Chandra. Scaling parameter-constrained language
- models with quality data. ArXiv, abs/2410.03083, 2024. URL https://api.semanticscholar.
- org/CorpusID:273162494.
- Jiefeng Chen, Timothy Nguyen, Dilan Gorur, and Arslan Chaudhry. Is forgetting less a good
- inductive bias for forward transfer? ArXiv, abs/2303.08207, 2023a. URL https://api.
- semanticscholar.org/CorpusId:257532608.
- Mayee F. Chen, Nicholas Roberts, K. Bhatia, Jue Wang, Ce Zhang, Frederic Sala, and Christopher Ré.
- Skill-it! a data-driven skills framework for understanding and training language models. ArXiv,
- abs/2307.14430, 2023b. URL https://api.semanticscholar.org/CorpusId:260203057.
- Ronen Eldan and Yuanzhi Li. Tinystories: How small can language models be and still speak coherent english?, 2023. URL https://arxiv.org/abs/2305.07759.
- Jean Gailly and Mark Adler. GNU gzip. GNU Operating System, 1992.
- 193 Claudio Greco, Barbara Plank, R. Fernández, and R. Bernardi. Psycholinguistics meets continual
- learning: Measuring catastrophic forgetting in visual question answering. In Annual Meeting of
- the Association for Computational Linguistics, 2019. URL https://api.semanticscholar.
- org/CorpusId:184488333.
- Jinghan He, Haiyun Guo, Kuan Zhu, Zihan Zhao, Ming Tang, and Jinqiao Wang. Seekr: Se-
- lective attention-guided knowledge retention for continual learning of large language mod-
- els. In Conference on Empirical Methods in Natural Language Processing, 2024. URL
- 200 https://api.semanticscholar.org/CorpusID:273901289.
- Heng Huang, Li Shen, Enneng Yang, and Zhenyi Wang. A comprehensive survey of forgetting in
- deep learning beyond continual learning. IEEE Transactions on Pattern Analysis and Machine
- 203 Intelligence, 47:1464-1483, 2023. URL https://api.semanticscholar.org/CorpusId:
- 204 259951356.
- Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, Stan-
- ley Jungkyu Choi, and Minjoon Seo. Towards continual knowledge learning of language mod-
- els. ArXiv, abs/2110.03215, 2021. URL https://api.semanticscholar.org/CorpusId:
- 208 238419458.
- Joel Jang, Seonghyeon Ye, Changho Lee, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun
- Kim, and Minjoon Seo. Temporalwiki: A lifelong benchmark for training and evaluating ever-
- evolving language models. In Conference on Empirical Methods in Natural Language Processing,
- 2022. URL https://www.aclanthology.org/2022.emnlp-main.418.pdf.
- 213 Khimya Khetarpal, M. Riemer, I. Rish, and Doina Precup. Towards continual reinforcement learning:
- A review and perspectives. J. Artif. Intell. Res., 75:1401-1476, 2020. URL https://api.
- semanticscholar.org/CorpusId:229679944.

- D. Kudithipudi, Mario Aguilar-Simon, Jonathan Babb, M. Bazhenov, Douglas Blackiston, J. Bongard, 216 Andrew P. Brna, Suraj Chakravarthi Raja, Nick Cheney, J. Clune, A. Daram, Stefano Fusi, Peter 217 Helfer, Leslie M. Kay, Nicholas A. Ketz, Z. Kira, Soheil Kolouri, J. Krichmar, Sam Kriegman, 218 Michael Levin, Sandeep Madireddy, Santosh Manicka, Ali Marjaninejad, Bruce L. McNaughton, 219 R. Miikkulainen, Zaneta Navratilova, Tej Pandit, Alice Parker, Praveen K. Pilly, S. Risi, T. Se-220 jnowski, Andrea Soltoggio, Nicholas Soures, A. Tolias, Darío Urbina-Meléndez, F. Valero-Cuevas, 221 Gido M. van de Ven, J. Vogelstein, Felix Wang, Ron Weiss, A. Yanguas-Gil, Xinyun Zou, and 222 H. Siegelmann. Biological underpinnings for lifelong learning machines. Nature Machine Intelli-223 gence, 4:196 - 210, 2022. URL https://doi.org/10.1038/s42256-022-00452-0. 224
- Jeffrey Li, Mohammadreza Armandpour, Iman Mirzadeh, Sachin Mehta, Vaishaal Shankar, Raviteja
 Vemulapalli, Samy Bengio, Oncel Tuzel, Mehrdad Farajtabar, Hadi Pouransari, and Fartash Faghri.
 Tic-lm: A web-scale benchmark for time-continual llm pretraining. ArXiv, abs/2504.02107, 2025.
 URL https://api.semanticscholar.org/CorpusId:277510618.
- Jia Liu, Jinguo Cheng, Xiangming Fang, Zhenyuan Ma, and Yuankai Wu. Evaluating temporal plasticity in foundation time series models for incremental fine-tuning. *ArXiv*, abs/2504.14677, 2025. URL https://api.semanticscholar.org/CorpusId:277954770.
- Hadi Nekoei, Akilesh Badrinaaraayanan, Aaron C. Courville, and Sarath Chandar. Continuous coordination as a realistic scenario for lifelong learning. *ArXiv*, abs/2103.03216, 2021. URL https://api.semanticscholar.org/CorpusId:232110854.
- Anastasia Razdaibiedina, Yuning Mao, Rui Hou, Madian Khabsa, Mike Lewis, and Amjad Almahairi.
 Progressive prompts: Continual learning for language models. *ArXiv*, abs/2301.12314, 2023. URL https://api.semanticscholar.org/CorpusID:256390383.
- Corban G. Rivera, C. Ashcraft, Alexander New, J. Schmidt, and Gautam K. Vallabha. Latent properties of lifelong learning systems. *ArXiv*, abs/2207.14378, 2022. URL https://api.semanticscholar.org/CorpusId:251196582.
- Samuel Schmidgall, Jascha Achterberg, Thomas Miconi, Louis Kirsch, Rojin Ziaei, S. P. Hajiseyedrazi, and Jason Eshraghian. Brain-inspired learning in artificial neural networks: a review. *ArXiv*, abs/2305.11252, 2023. URL https://api.semanticscholar.org/CorpusId:258823273.
- Thomas Scialom, Tuhin Chakrabarty, and Smaranda Muresan. Fine-tuned language models are continual learners. In *Conference on Empirical Methods in Natural Language Processing*, 2022. URL https://api.semanticscholar.org/CorpusID:252815378.
- Haizhou Shi, Zihao Xu, Hengyi Wang, Weiyi Qin, Wenyuan Wang, Yibin Wang, and Hao Wang.
 Continual learning of large language models: A comprehensive survey. *ACM Computing Surveys*,
 2024. URL https://api.semanticscholar.org/CorpusId:269362836.
- U.S. Department of Health & Human Services, Administration for Children & Families. Head start early learning outcomes framework: Ages birth to five, 2024. URL https://headstart.gov/school-readiness/article/head-start-early-learning-outcomes-framework. Last updated: December 23, 2024; Accessed: 2025-08-20.
- Xiao Wang, Yuan Zhang, Tianze Chen, Songyang Gao, Senjie Jin, Xianjun Yang, Zhiheng Xi,
 Rui Zheng, Yicheng Zou, Tao Gui, Qi Zhang, and Xuanjing Huang. Trace: A comprehensive
 benchmark for continual learning in large language models. ArXiv, abs/2310.06762, 2023. URL
 https://api.semanticscholar.org/CorpusID:263830425.
- Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan, Thuy-Trang Vu, and Gholamreza Haffari.
 Continual learning for large language models: A survey. *ArXiv*, abs/2402.01364, 2024. URL https://api.semanticscholar.org/CorpusId:267406164.
- Yuhao Wu, Tongjun Shi, Karthick Sharma, Chun Seah, and Shuhao Zhang. Online continual knowledge learning for language models. *ArXiv*, abs/2311.09632, 2023. URL https://api.semanticscholar.org/CorpusId:265221422.
- Yongxin Xu, Philip S. Yu, Zexin Lu, Xu Chu, Yujie Feng, Bo Liu, and Xiao-Ming Wu. Klf: Knowledge localization and fusion for language model continual learning. 2024. URL https://api.semanticscholar.org/CorpusId:271843361.

- Lixiang Yan, Samuel Greiff, Ziwen Teuber, and D. Gaević. Promises and challenges of generative artificial intelligence for human learning. *Nature human behaviour*, 8 10:1839–1850, 2024. URL https://api.semanticscholar.org/CorpusId:271924303.
- Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. Character-level convolutional networks for
 text classification. In *Neural Information Processing Systems*, 2015. URL https://api.semanticscholar.org/CorpusID:368182.
- Hongbo Zhao, Fei Zhu, Rundong Wang, Gaofeng Meng, and Zhaoxiang Zhang. Mllm-cl: Continual learning for multimodal large language models. *ArXiv*, abs/2506.05453, 2025. URL https://api.semanticscholar.org/CorpusId:279243888.
- Junhao Zheng, Xidi Cai, Qiuke Li, Duzhen Zhang, Zhongzhi Li, Yingying Zhang, Le Song, and
 Qianli Ma. Lifelongagentbench: Evaluating llm agents as lifelong learners. ArXiv, abs/2505.11942,
 2025. URL https://api.semanticscholar.org/CorpusId:278739762.

279 A Dataset Construction

Figure 2 gives an overview of our dataset, including examples of skills, subskills, goal and indicator.
We also present an example of an edge from the skill graph. Figure 3 shows the number of incoming and outgoing edges from each stage. Figure 4 explains the data construction process and gives examples from each stage of the data generation pipeline.

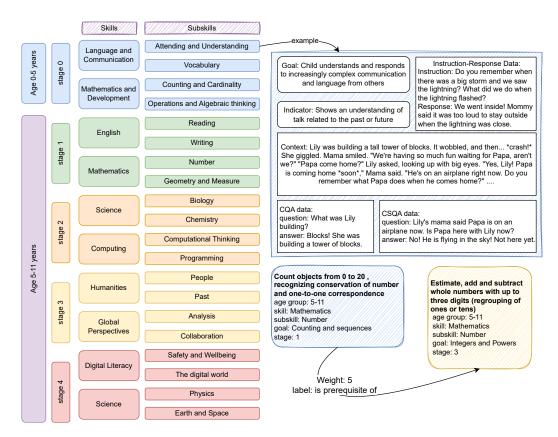


Figure 2: Developmental framework for children aged 0-11 years, categorized into stages (0-4). Only examples of skills and subskills are mentioned here. An example of how the data looks like is given in the top right. Two nodes and an edge from the skill graph is given in the bottom right.

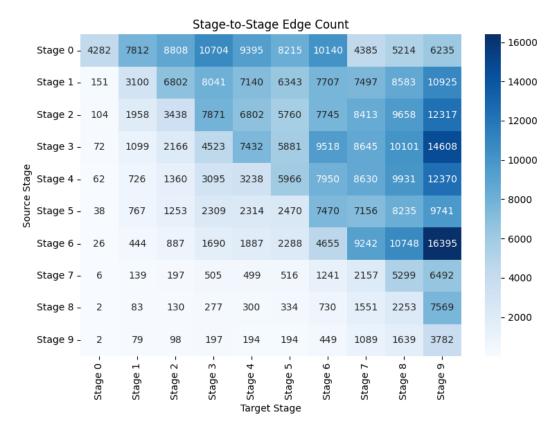


Figure 3: Heatmap showing the number of prerequisite edges between stages in the skill graph. Rows correspond to source stages, columns to target stages, and color intensity indicates the number of connections.

284 B Data Verification

For both the methods, 500K texts are sampled from each of the Paragraphs and Instruction-response pairs.

287 B.1 Diversity

For the diversity of the text, we follow Chang et al. [2024] and calculate the compression ratio of the text as

$$\mathrm{CR}(D) = \frac{Original size of D\ (bytes)}{Compressed size of D\ (bytes)},$$

290 and define diversity by

291

292

293

294

295

$$Dr(D) = 1/CR(D).$$

A higher compression ratio $\mathrm{CR}(D)$ indicates greater redundancy, meaning lower diversity in the text. Thus, diversity $\mathrm{Dr}(D)$ increases when redundancy decreases. We see diversity ranging between 30.77% and 35.60%, which is similar to other work. As a comparison, we also calculated the diversity of 500K samples from the validation set of TinyStories, a paper exploring synthetic data generation to train a small language model. Their text diversity ranges from 31.04% to 32.66% within the pretraining and instruct data, respectively Eldan and Li [2023].

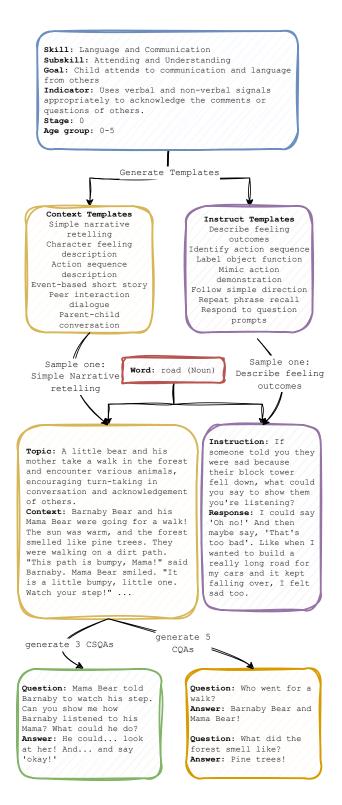


Figure 4: Synthetic data generation pipeline

Table 2: Diversity and Deduplication metrics for context and instruction-response data across stages

| Stage | Con | ntext | IR | | |
|-------|-------------|--------|--------|-----------------------------|--|
| | Div ↑ Dedup | | Div ↑ | $\mathbf{Dedup} \downarrow$ | |
| 0 | 34.29% | 11.83% | 30.77% | 3.50% | |
| 1 | 35.60% | 5.36% | 31.73% | 3.85% | |
| 2 | 34.17% | 15.47% | 32.64% | 2.54% | |
| 3 | 34.68% | 14.86% | 32.97% | 2.09% | |
| 4 | 35.45% | 13.41% | 33.14% | 1.93% | |

B.2 Deduplication

297

For semantic deduplication⁶, we pass the texts through a sentence encoder and find the deduplication rate as the percentage of sentences that have cosine similarity of at least 0.95 with another sentence in the same stage.

Table 3: Diversity and Deduplication Rates when Considering Pairwise Stages

| Stage Pair | Co | ntext |
|------------|--------|-----------------------------|
| | Div ↑ | $\mathbf{Dedup} \downarrow$ |
| | | |
| 0, 1 | 31.29% | 0.3% |
| 0, 2 | 31.96% | 0.1% |
| 0, 3 | 32.25% | 0.0% |
| 0, 4 | 32.50% | 0.0% |
| 1, 2 | 32.27% | 0.3% |
| 1, 3 | 32.52% | 0.2% |
| 1, 4 | 32.71% | 0.1% |
| 2, 3 | 32.82% | 0.4% |
| 2, 4 | 32.94% | 0.2% |
| 3, 4 | 33.07% | 0.2% |

B.3 Detailed Readability Metrics

Note that average grade of the data is slightly higher than the intended age of the data (especially for the first few stages). However, this is because not all skills we generate data for are, in real-life, text-based. Thus, demonstrating them in language ends up requiring complex words, which affects the readability score. For example, children can verbally reason about cause-and-effect in multi-turn conversations, but when written down, that same dialogue is rated at a much higher reading level than the child can actually read, leading to higher readability scores in our data.

Table 4: Average readability scores of generated data across stages, reported for context, comprehension QA (CQA), skill-testing QA (CSQA), and instruction–response (IR) data. Scores generally increase with stage, reflecting controlled growth in textual complexity aligned with developmental progression

| Stage | Context | CQA | CSQA | IR |
|-------|-----------|-----------|-----------|-----------|
| 0 | 4.61 1.87 | 2.38 2.88 | 3.07 2.26 | 4.48 1.52 |
| 1 | 5.24 1.72 | 4.39 1.81 | 4.44 1.62 | 4.86 1.41 |
| 2 | 5.18 1.93 | 4.39 1.80 | 4.69 1.54 | 4.69 1.59 |
| 3 | 5.51 1.85 | 4.65 1.70 | 4.98 1.46 | 5.03 1.50 |
| 4 | 6.42 1.79 | 5.63 1.44 | 5.96 1.30 | 5.91 1.34 |

307

⁶We use the following repo for semantic deduplication: https://github.com/MinishLab/semhash

Table 5: Detailed Readability Metrics Across all 5 Stages and Datasets

| Dataset | Stage | Flesch Kincaid | SMOG | Coleman Liau | Automated Readability | Dale Chall | Gunning Fog |
|---------|-------|----------------|-----------|--------------|-----------------------|------------|-------------|
| Context | 0 | 3.15 0.35 | 6.90 0.76 | 4.28 0.51 | 1.68 0.47 | 6.69 0.27 | 4.94 0.34 |
| Context | 1 | 3.68 0.35 | 7.55 0.73 | 5.18 0.50 | 2.58 0.49 | 6.70 0.29 | 5.74 0.35 |
| Context | 2 | 3.80 0.36 | 7.54 0.75 | 4.21 0.46 | 2.25 0.48 | 7.18 0.35 | 6.12 0.38 |
| Context | 3 | 4.16 0.36 | 7.84 0.74 | 4.58 0.48 | 2.71 0.50 | 7.27 0.36 | 6.48 0.38 |
| Context | 4 | 5.13 0.42 | 8.76 0.79 | 5.39 0.51 | 3.77 0.56 | 7.89 0.34 | 7.58 0.45 |
| CQA | 0 | 0.79 0.35 | 5.06 0.59 | 0.26 0.59 | -1.47 0.43 | 6.89 0.30 | 2.75 0.34 |
| CQA | 1 | 2.73 0.37 | 6.45 0.59 | 4.10 0.54 | 1.65 0.49 | 6.47 0.26 | 4.92 0.45 |
| CQA | 2 | 2.74 0.38 | 6.42 0.59 | 4.00 0.53 | 1.67 0.50 | 6.44 0.28 | 5.07 0.45 |
| CQA | 3 | 3.04 0.37 | 6.66 0.59 | 4.37 0.52 | 2.08 0.49 | 6.41 0.27 | 5.36 0.44 |
| CQA | 4 | 4.08 0.38 | 7.54 0.59 | 5.59 0.48 | 3.52 0.49 | 6.50 0.25 | 6.54 0.47 |
| CSQA | 0 | 1.34 0.28 | 5.37 0.70 | 2.07 0.43 | -0.20 0.36 | 6.21 0.20 | 3.65 0.27 |
| CSQA | 1 | 2.84 0.30 | 6.36 0.71 | 4.14 0.37 | 2.04 0.40 | 6.03 0.21 | 5.24 0.33 |
| CSQA | 2 | 3.16 0.29 | 6.54 0.72 | 4.33 0.37 | 2.43 0.39 | 6.08 0.23 | 5.59 0.33 |
| CSQA | 3 | 3.49 0.29 | 6.81 0.70 | 4.64 0.37 | 2.87 0.39 | 6.14 0.25 | 5.96 0.32 |
| CSQA | 4 | 4.62 0.33 | 7.72 0.72 | 5.56 0.41 | 4.25 0.46 | 6.50 0.27 | 7.12 0.37 |
| IR | 0 | 2.97 0.47 | 6.32 0.64 | 4.12 0.52 | 2.25 0.62 | 5.81 0.23 | 5.43 0.48 |
| IR | 1 | 3.40 0.45 | 6.61 0.65 | 4.51 0.50 | 2.88 0.61 | 5.76 0.25 | 6.02 0.50 |
| IR | 2 | 3.16 0.37 | 6.62 0.72 | 4.23 0.45 | 2.33 0.50 | 6.10 0.26 | 5.68 0.43 |
| IR | 3 | 3.55 0.37 | 6.93 0.71 | 4.62 0.46 | 2.87 0.51 | 6.13 0.27 | 6.09 0.42 |
| IR | 4 | 4.59 0.41 | 7.66 0.73 | 5.41 0.46 | 4.13 0.56 | 6.46 0.28 | 7.20 0.47 |

308 C Hyperparameters

All experiments were conducted with a consistent set of training hyperparameters to ensure com-309 parability across runs. Models were initialized using the kaiming normal method unless otherwise 310 specified, and trained with AdamW optimizer ($\beta_1 = 0.9, \beta_2 = 0.98, \epsilon = 1e - 8$) with weight decay 311 of 0.01. We used a base learning rate of 5e-3, applied gradient clipping with a maximum norm 312 of 1.0. We used gradient accumulation (8 steps with batch size 24 on 8 GPUs, yielding an effective 313 batch size of 1536). Training was performed for one full epoch over each dataset split with a context 314 length of 1024 tokens. Mixed precision was enabled with bfloat16 (bf16) for efficiency, while fp16 315 was disabled. All experiments were seeded with 42 for reproducibility. For inference, the model was 316 loaded in bfloat16 precision with padding set to the EOS token and leftside padding for alignment. 317 Prompts were tokenized with a maximum length of 512 tokens, and generation used a temperature of 0.7, top-p sampling of 0.95, and a maximum of 128 new tokens per prompt.

320 D Results

Table 6 gives the results of all experiments on IR test set. Table 7 gives the results of all experiments on CQA test set. Table 8 gives the results of all experiments on CSQA test set. Per-stage per-Indicator results can be found here: Results sheet. Forgetting analysis is shown in Figure 5. Relation of forgetting analysis to the skill graph can be drawn from Figure 6.

E Prompts

325

326

328

330

331 332

E.1 Edge Prediction

327 System prompt for Edge prediction

You are an expert in skill development and cognitive science. Your task is to analyze the relationship between two skill indicators and determine if there is a logical prerequisite dependency between them.

Each skill indicator is given with:

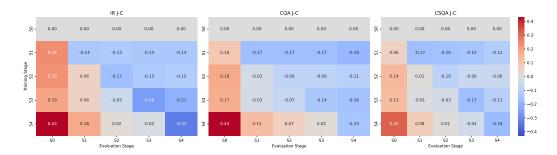


Figure 5: Forgetting analysis across training setups. The plots show performance differences between joint and continual training for IR, CQA, and CSQA test sets across stages 0–4. The Y-axis corresponds to models trained upto a stage. The X-axis corresponds to test set of mentioned stage.

Table 6: All results for IR test set. The column represents each stage on which a model is being evaluated.

| Test type | IR (rating out of 5) | | | | | |
|-----------------|----------------------|------|------|------|------|--|
| Stages | 0 | 1 | 2 | 3 | 4 | |
| M_0 | 4.16 | 3.29 | 2.97 | 2.83 | 2.49 | |
| M_1 | 3.70 | 3.70 | 3.21 | 3.08 | 2.80 | |
| M_2 | 3.71 | 3.55 | 3.56 | 3.27 | 3.00 | |
| M_3 | 3.64 | 3.45 | 3.35 | 3.57 | 3.07 | |
| M_4 | 3.38 | 3.35 | 3.32 | 3.34 | 3.55 | |
| M_{012} | 4.22 | 3.81 | 3.55 | 3.34 | 3.07 | |
| M_{01} | 4.19 | 3.73 | 3.25 | 3.12 | 2.84 | |
| M_{0-1} | 3.94 | 3.87 | 3.38 | 3.26 | 2.98 | |
| M_{0123} | 4.15 | 3.79 | 3.56 | 3.55 | 3.14 | |
| M_{0-1-2} | 3.99 | 3.75 | 3.72 | 3.47 | 3.19 | |
| M_{01234} | 4.16 | 3.80 | 3.60 | 3.60 | 3.46 | |
| $M_{0-1-2-3}$ | 3.97 | 3.73 | 3.61 | 3.82 | 3.34 | |
| $M_{0-1-2-3-4}$ | 3.73 | 3.63 | 3.58 | 3.62 | 3.78 | |

Table 7: All results for CQA test set. The column represents each stage on which a model is being evaluated.

| Test type | | CQA (1 | rating o | ut of 5) | |
|-----------------|------|--------|----------|----------|------|
| Stages | 0 | 1 | 2 | 3 | 4 |
| M_0 | 4.16 | 3.29 | 2.97 | 2.83 | 2.49 |
| M_1 | 3.70 | 3.70 | 3.21 | 3.08 | 2.80 |
| M_2 | 3.71 | 3.55 | 3.56 | 3.27 | 3.00 |
| M_3 | 3.64 | 3.45 | 3.35 | 3.57 | 3.07 |
| M_4 | 3.38 | 3.35 | 3.32 | 3.34 | 3.55 |
| M_{012} | 4.22 | 3.81 | 3.55 | 3.34 | 3.07 |
| M_{01} | 4.19 | 3.73 | 3.25 | 3.12 | 2.84 |
| M_{0-1} | 3.94 | 3.87 | 3.38 | 3.26 | 2.98 |
| M_{0123} | 4.15 | 3.79 | 3.56 | 3.55 | 3.14 |
| M_{0-1-2} | 3.99 | 3.75 | 3.72 | 3.47 | 3.19 |
| M_{01234} | 4.61 | 4.27 | 4.05 | 3.87 | 3.45 |
| $M_{0-1-2-3}$ | 4.42 | 4.27 | 4.09 | 3.97 | 3.45 |
| $M_{0-1-2-3-4}$ | 4.17 | 4.14 | 3.97 | 3.85 | 3.60 |

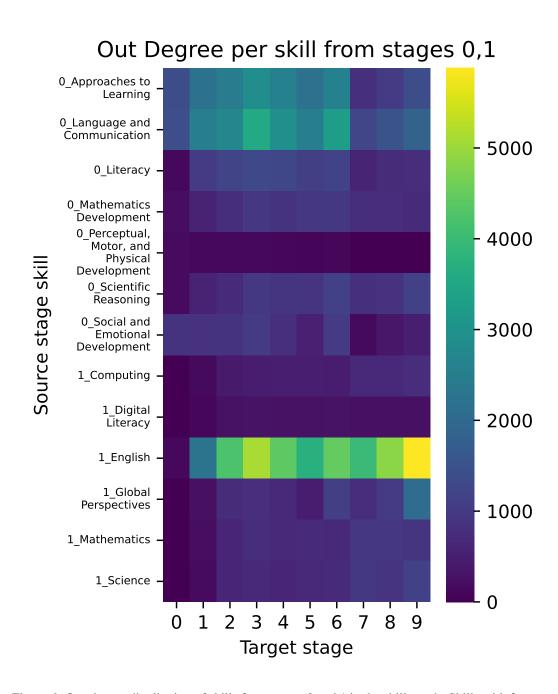


Figure 6: Out-degree distribution of skills from stages 0 and 1 in the skill graph. Skills with fewer outgoing prerequisite edges (e.g., Perceptual, Motor, and Physical Development; Digital Literacy) are less connected to later stages and are observed to be more vulnerable to forgetting in continual training.

Table 8: All results for CSQA test set. The column represents each stage on which a model is being evaluated.

| Test type | (| CSQA (| rating o | out of 5 |) |
|-----------------|------|--------|----------|----------|------|
| Stages | 0 | 1 | 2 | 3 | 4 |
| M_0 | 3.89 | 2.85 | 2.52 | 2.33 | 1.95 |
| M_1 | 3.63 | 3.35 | 2.92 | 2.75 | 2.39 |
| M_2 | 3.53 | 3.25 | 3.15 | 2.87 | 2.53 |
| M_3 | 3.51 | 3.22 | 3.03 | 3.10 | 2.61 |
| M_4 | 3.29 | 3.13 | 3.00 | 2.93 | 2.89 |
| M_{012} | 3.97 | 3.48 | 3.21 | 2.96 | 2.61 |
| M_{01} | 3.93 | 3.37 | 2.93 | 2.76 | 2.40 |
| M_{0-1} | 3.87 | 3.55 | 3.09 | 2.91 | 2.51 |
| M_{0123} | 3.97 | 3.47 | 3.21 | 3.08 | 2.65 |
| M_{0-1-2} | 3.83 | 3.47 | 3.31 | 3.03 | 2.66 |
| M_{01234} | 3.97 | 3.49 | 3.24 | 3.13 | 2.88 |
| $M_{0-1-2-3}$ | 3.83 | 3.48 | 3.24 | 3.26 | 2.76 |
| $M_{0-1-2-3-4}$ | 3.65 | 3.41 | 3.23 | 3.17 | 3.05 |

```
- a_label and a_id
    - b_label and b_id
335
336
    These represent two distinct skill indicators. You must determine whether one is a
337
338
        prerequisite for the other.
339
    Instructions:
340
    - A skill X is a prerequisite for skill Y if Y logically requires understanding or
341
342
         demonstrating X beforehand.
    - Compare the meaning of a_label and b_label to determine if:
343
      - A depends on B edge from b_id to a_id
344
      - B depends on A edge from a_id to b_id
345
      - No clear dependency no edge
346
347
    Output format:
348
    Return a JSON object like:
349
350
    ""json
351
    }}
352
      "edge": true or false,
353
      "from": "source_id" or "NA",
354
      "to": "target_id" or "NA",
355
      "reason": "Brief explanation of the dependency or lack thereof"
356
357
358
359
    - If there is a dependency, set edge: true, from as the prerequisite's ID, and to as
360
          the dependent's ID.
361
    - If there is no clear prerequisite relationship, set edge: false and "from": "NA",
362
         "to": "NA" with a brief justification in reason.
363
364
365
    Only base your answer on the textual meaning of the labels, and only report direct
         dependencies (not transitive or indirect ones).
369
```

User prompt for Edge prediction

```
Given the following skill indicators:
- a_label: {label_1}
- a_id: {id_1}
- b_label: {label_2}
- b_id: {id_2}
```

```
Determine the dependency relationship and output the JSON:
376
377
    "'json
378
    }}
379
      "edge": true or false,
380
      "from": "source_id" or "NA",
381
      "to": "target_id" or "NA",
382
      "reason": "Brief explanation of the dependency or lack thereof"
383
    }}
384
385
```

387 E.2 Edge weight prediction

System prompt:

388

```
You are an expert in child development, skill acquisition, and cognitive science.
390
         Your task is to rate the strength of a prerequisite relationship between two
391
         skill indicators. Each input includes:
392
    - from_label and to_label: the skill indicators (already determined to be in a
393
        prerequisite relationship, where from_label is a prerequisite for to_label)
394
395
    - Additional metadata: age groups, subskills, goals, developmental stages, and a
        rationale for why the edge exists.
396
397
398
    Instructions:
    Rate the dependency strength on a scale from 1 to 5, where:
399
    - 1 = Very weak dependency (minimal or contextual support, can often be developed
400
401
         independently)
    - 2 = Weak dependency (some support role, but not always required)
402
    - 3 = Moderate dependency (often occurs first, but not strictly necessary)
403
404
    - 4 = Strong dependency (usually needed before progressing)
    - 5 = Very strong dependency (essential foundational step for the next)
405
406
    Your response should consider:
407
    1. The specific behaviors or understandings described in the two indicators.
408
    2. Whether the earlier skill is conceptually or procedurally required to perform the
409
410
    3. The closeness of developmental stages and subskills.
411
412
    Output Format:
413
    Return your decision as a JSON object:
414
    ""json
415
    }}
416
      "weight": [an integer from 1 to 5],
417
      "reason": "[a brief explanation of why this weight reflects the strength of the
418
          dependency] "
419
    }}
420
    ""
421
```

User prompt:

```
424
    Given the following information about a prerequisite relationship between two skill
         indicators:
426
427
428
     - from_label: {from_label}
     - from_id: {from_id}
429
        - age group: {from_age_group}
430
431
        - skill: {from_skill}
        - subskill: {from_subskill}
432
        - goal: {from_goal}
433
434
        - stage: {from_stage}
435
436
437
    - to_label: {to_label}
```

```
- to_id: {to_id}
439
        - age group: {to_age_group}
440
        - skill: {to_skill}
441
        - subskill: {to_subskill}
442
        - goal: {to_goal}
443
        - stage: {to_stage}
444
445
    This relationship has already been labeled as a prerequisite edge (from_id to_id).
446
447
    Rationale for this dependency:
448
449
    "{reason}"
450
    Rate the strength of this dependency on a scale from 1 to 5.
451
452
    Output a JSON object:
453
    "'json
454
    {{
455
      "weight": [an integer from 1 to 5],
456
      "reason": "Brief explanation of why this weight reflects the strength of the
457
458
           dependency"
    }}
459
    ""
469
```

E.3 Templates

463

System prompt for generating templates for IR data:

```
You are an expert in child development, skill acquisition, curriculum design, and
465
         language model pretraining. Your task is to identify developmentally
466
         appropriate and general **non-instructional text types** for synthetic
467
468
        pretraining of a language model.
469
470
    Each input includes:
471
    - indicator: a natural language description of the learning objective or task
    - age_group: developmental age (e.g., 05, 511, 1114)
472
473
    - skill: broad academic or developmental domain (e.g., Mathematics, English,
474
         Scientific Reasoning)
    - subskill: a specific subdomain or area of focus (e.g., Listening, Measurement,
475
        Problem-solving)
476
    - goal: the purpose or nature of the learning (e.g., Application, Reflection,
477
478
479
    - stage: the curriculum stage (0 to 9, loosely corresponding to increasing age and
480
        complexity)
481
    Instructions:
482
    Return a list of **general non-instructional text types** that:
483
484
    - Are suitable for the learner's developmental stage
    - Reflect naturalistic or structured formats that don't rely on explicit
485
         instructionresponse pairs
486
    - Can be used as abstract templates to generate content across many topics
487
    - Are defined at a high level of abstraction (e.g., "peer dialogue", "narrative
488
         description", "cause-effect explanation")
489
490
491
    **CRITICALLY IMPORTANT**:
    - Provide format categories, NOT specific content or scenarios
492
    - Text types should be 2-5 words that describe a general format, not complete
493
         sentences
494
    - Each text type should be usable with ANY topic relevant to the age/skill
495
496
497
    **Examples of appropriate non-instructional text types**:
498
    - "Narrative story with characters"
499
500
    - "Peer conversation transcript"
    - "Process description passage"
```

```
- "Personal reflection monologue"
502
    **Examples of inappropriate text types** (too specific):
504
    - "Story about a child going to the zoo"
505
    - "Conversation between friends about toys"
506
    - "Description of a butterfly's life cycle"
507
508
    Output Format:
509
    Return your result as a JSON object with the following structure:
510
511
    "'json
512
    }}
513
      "text_types": ["...", "...", "..."]
514
    }}
515
516
517
    Ensure the list is:
518
    - 1520 items long
519
    - Abstract enough to work across many topics
520
    - Varied across narration, description, interaction, emotion, reasoning
521
    - Appropriate in complexity for the given age group and learning goal
522
523
    Only output the JSON object.
524
```

User prompt for generating templates for IR data:

```
Given the following information about a learning objective, return a list of general
528
         , reusable non-instructional text formats that can serve as templates for
529
         synthetic training data:
530
    - indicator: {indicator}
532
    - age_group: {age_group}
533
    - skill: {skill}
534
    - subskill: {subskill}
    - goal: {goal}
    - stage: {stage}
537
538
    IMPORTANT: Provide ABSTRACT FORMAT CATEGORIES (2-5 words each), not specific content
539
          or scenarios.
540
541
    Examples of good non-instructional formats:
542
    - "Peer dialogue transcript"
543
544
    - "Sequential process description"
    - "Character-driven narrative"
545
    - "Emotional experience monologue"
546
547
    Examples of unsuitable formats (too specific):
548
    - "Conversation between friends about toys"
    - "Description of a butterfly's life cycle"
550
    - "Story about going to the beach"
551
552
    Ensure your list contains:
    - 15 to 20 developmentally appropriate text formats
554
555
    - General templates that can be combined with ANY relevant topic
    - Varied format types that don't rely on explicit instruction-response pairs
556
557
    Return only a JSON object in the following format:
559
    ""json
560
561
      "text_types": ["...", "...", "..."]
    }}
563
    ""
564
```

566 System prompt for generating templates for Context data:

```
567
568
    You are an expert in child development, skill acquisition, curriculum design, and
         language model pretraining. Your task is to identify developmentally
569
         appropriate and general **instruction-response text types** for synthetic
570
571
        pretraining of a language model.
572
    Each input includes:
573
    - indicator: a natural language description of the learning objective or task
574
    - age_group: developmental age (e.g., 05, 511, 1114)
575
      skill: broad academic or developmental domain (e.g., Mathematics, English,
576
577
         Scientific Reasoning)
    - subskill: a specific subdomain or area of focus (e.g., Listening, Measurement,
578
        Problem-solving)
579
580
    - goal: the purpose or nature of the learning (e.g., Application, Reflection,
581
         Evaluation)
    - stage: the curriculum stage (0 to 9, loosely corresponding to increasing age and
582
         complexity)
583
584
    Instructions:
585
    Return a list of **general instruction-response style text types** that:
586
    - Are suitable for the learner's developmental stage
587
     - Can be used in instruction tuning and task-based language modeling
588
589
    - Involve a clearly defined instruction format that can be applied across many
590
         topics
    - Are defined at a high level of abstraction (e.g., "explain why X occurs", "compare
591
          and contrast X and Y")
592
593
594
    **CRITICALLY IMPORTANT**:
    - Provide abstract instruction formats, NOT specific prompts or questions
595
    - Text types should be 2-5 words describing a general instruction format
596
597
    - Each text type should be usable with ANY topic relevant to the age/skill
         combination
598
599
    **Examples of appropriate instruction-response text types**:
600
    - "Compare and contrast analysis"
601
    - "Explain why reasoning"
603
    - "Step-by-step instruction"
    - "Open-ended reflection prompt"
604
605
606
    **Examples of inappropriate text types** (too specific):
607

    "Explain why plants need water"

    - "Compare dogs and cats"
608
    - "Describe your favorite toy"
609
610
611
612
    Return your result as a JSON object with the following structure:
613
    ""json
614
615
    {{
      "text_types": ["...", "...", "..."]
616
    }}
617
    "
618
619
    Ensure the list is:
621
    - 1520 items long
    - Abstract enough to work across many topics
622
    - Varied across explanation, reasoning, reflection, comparison, instruction,
623
         imagination
    - Appropriate in complexity for the given age group and learning goal
625
626
    Only output the JSON object.
627
```

9 User prompt for generating templates for Context data:

```
Given the following information about a learning objective, return a list of general
632
         , reusable instruction-response text formats that can serve as templates for
         synthetic training data:
633
634
635
    - indicator: {indicator}
    - age_group: {age_group}
636
    - skill: {skill}
637
638
    - subskill: {subskill}
    - goal: {goal}
639
640
    - stage: {stage}
641
    IMPORTANT: Provide ABSTRACT INSTRUCTION FORMATS (2-5 words each), not specific
642
         questions or prompts.
643
    Examples of good instruction formats:
645
    - "Compare and contrast analysis"
646
647
    - "Explain why reasoning"
    - "Problem-solving walkthrough"
648
    - "Open-ended reflection prompt"
649
650
    Examples of unsuitable formats (too specific):
651
652
     "Explain why plants need water"
    - "Compare dogs and cats"
653
    - "Solve this math problem"
654
655
    Ensure your list contains:
656
    - 15 to 20 developmentally appropriate instruction formats
657
658
    - General templates that can be combined with ANY relevant topic
    - Varied instruction types that address different cognitive processes
659
660
661
    Return only a JSON object in the following format:
662
    ""json
663
664
      "text_types": ["...", "...", "..."]
665
666
668
```

E.4 Context

670

System prompt for generating context data:

```
You are an AI model generating training data to help language models simulate human
672
         developmental skills at various stages from early childhood through early
673
        adolescence.
674
675
    Your task is to create engaging, developmentally appropriate texts based on provided
676
          developmental indicators, skills, and a tuple of word and its part of speech.
677
678
    Strictly follow these guidelines:
679
680
681
    1. **Developmental Appropriateness:**
       - Stage 0 (Age 5): Use simple sentences, concrete concepts, familiar experiences,
682
            present tense focus
683
684
       - Stages 1-3 (Ages 6-8): Introduce basic past/future concepts, simple cause-
685
           effect, familiar settings
       - Stages 4-6 (Ages 9-11): Include more complex reasoning, abstract thinking,
686
687
           varied sentence structures
       - Stages 7-9 (Ages 12-14): Incorporate hypothetical scenarios, multiple
688
           perspectives, sophisticated vocabulary
689
    2. **Context Generation:**
```

```
- Use the provided word and its part of speech to create a meaningful,
692
           developmentally appropriate topic
693
       - **Ensure the selected word and expanded topic fit the required Text Type
694
           Template (context_template)**
695
       - Expand the selected word into a more detailed, skill-aligned topic that
696
           resonates with the target age group
697
       - Generate a rich, complete, and engaging text matching the provided context
698
699
           template
       - The generated text must be **between 250 and 500 words regardless of
700
701
            developmental stage**
702
       - The text must clearly align with the skill, subskill, goal, and indicator
       - The selected word does not need to explicitly appear in the final text
703
704
705
    3. **Writing Style by Stage:**
       - **Early Stages (0-3):** Simple vocabulary, short to medium sentences, concrete
706
            experiences, repetitive patterns for reinforcement
707
       - **Middle Stages (4-6):** More varied vocabulary, complex sentences,
708
           introduction of abstract concepts, problem-solving scenarios
709
       - **Later Stages (7-9): ** Sophisticated vocabulary, complex sentence structures,
710
            abstract reasoning, multiple viewpoints
711
712
    4. **Content Enrichment:**
713
714
       - Include age-appropriate actions, feelings, interactions, and sensory details
       - Incorporate social situations relevant to the developmental stage
715
716
       - Use scenarios that promote the specific skill being targeted
       - Avoid overly abstract or culturally specific references unless appropriate for
717
           the age group
718
719
    5. **Output Format: ** Strictly return the output in the following JSON structure:
720
    ""json
721
    {{
722
        "expanded_topic": "<expanded topic>",
723
        "generated_text": "<generated text between 250 and 500 words>"
724
725
    }}
726
    Only output the JSON. No additional commentary.
727
```

User prompt for generating context data:

```
Generate a rich and engaging context text based on the following input:
731
732
    - ID: {id}
733
734
    - Indicator: {indicator}
    - Skill: {skill}
735
    - Sub-skill: {subskill}
736
    - Goal: {goal}
737
    - Age Group: {age_group}
738
739
    - Stage: {stage}
    - Text Type Template: {context_template}
740
    - (Word, Part of speech): {word_list}
741
742
743
    - Consider the developmental stage ({stage}) and age group ({age_group}) when
744
         crafting vocabulary, sentence complexity, and content themes
745
746
    - Expand the selected word into a skill-relevant topic **that fits the Text Type
747
    - Generate a detailed text of **250500 words** following the context template
748
    - Enrich the text with developmentally appropriate actions, emotions, and
749
        interactions
750
    - Ensure the content promotes the specific skill and subskill being targeted
751
752
    Output strictly in this format:
753
    ""json
754
755
    {{
        "expanded_topic": "<expanded topic>",
756
```

```
"generated_text": "<generated text between 250 and 500 words>"
| }}
| '''
```

```
E.5 CQA
761
    System prompt for generating CQA data:
762
    You are an AI model generating training data to help language models simulate human
765
         reading comprehension skills at various stages from early childhood through
         early adolescence.
766
767
    Your task is to create 5 developmentally appropriate question-answer pairs based on
768
         a provided text, ensuring all questions test understanding of the given
769
        paragraph and can be answered directly from the text.
770
771
    Strictly follow these guidelines:
772
773
774
    1. **Developmental Appropriateness by Stage:**
       - Stage 0 (Age 5): Simple "what/who/where" questions, literal comprehension,
775
            single-step reasoning
776
       - Stages 1-3 (Ages 6-8): Basic "why/how" questions, simple cause-effect, sequence
777
            understanding, character feelings
778
779
       - Stages 4-6 (Ages 9-11): Inference questions, comparing/contrasting, predicting
           outcomes, understanding motivations
780
       - Stages 7-9 (Ages 12-14): Complex analysis, multiple perspectives, abstract
781
782
            concepts, theme identification
783
    2. **Question Creation Standards:**
784
785
       - **All answers must be directly supported by information in the provided text**
       - No questions requiring outside knowledge or information not present in the text
786
       - Questions should test different types of comprehension appropriate to the
787
788
            developmental stage
       - Vary question types to assess different reading skills (literal, inferential,
789
790
            evaluative)
791
       - Use vocabulary and sentence complexity appropriate to the age group
       - Ensure questions are engaging and relevant to the child's interests and
792
            experiences
793
794
    3. **Question Types by Stage: **
795
796
       - **Early Stages (0-3):** Literal recall, identifying main characters/objects,
797
           simple sequence, basic emotions
       - **Middle Stages (4-6):** Cause-effect relationships, character motivations,
798
           comparing details, simple predictions
799
        **Later Stages (7-9): ** Drawing conclusions, analyzing relationships,
800
           evaluating actions, understanding themes
801
802
    4. **Answer Generation:**
803
       - Create authentic child responses that demonstrate comprehension at the target
804
            developmental stage
805
       - Use vocabulary and sentence structures appropriate to the age group
806
       - Include natural speech patterns and expressions typical of the developmental
807
808
       - Ensure answers are complete but not overly elaborate for the age group
809
       - Answers should sound conversational and natural, not textbook-like
810
811
    5. **Content Guidelines:**
812
813
```

- **Purely verbal exchanges** no references to physical gestures or non-verbal
 actions
- No formatting (bold, italics, markdown)

814

815

816 817

- Questions should flow naturally and cover different aspects of the text
- Ensure logical progression from simpler to more complex questions when appropriate

```
- Include a mix of question types (factual, inferential, personal connection when
819
             text-supported)
820
821
    6. **Quality Standards:**
822
       - Every question must be answerable using only information provided in the text
823
       - Questions should test genuine comprehension, not just memory of isolated facts
824
       - Avoid questions with obvious or trivial answers
825
       - Ensure questions are meaningful and help assess understanding of key text
826
            elements
827
828
       - Create questions that feel natural in an educational setting
829
    7. **Output Format: ** Strictly return the output in the following JSON structure:
830
    ""json
831
    {{
832
         "question_answer_pairs": [
833
834
                "question": "<question 1>",
835
                "answer": "<answer 1>"
836
            }},
837
            {{
838
                "question": "<question 2>",
839
                "answer": "<answer 2>"
840
            }},
841
842
                "question": "<question 3>",
843
                "answer": "<answer 3>"
844
            }},
845
846
                "question": "<question 4>",
847
                "answer": "<answer 4>"
848
            }},
849
850
                "question": "<question 5>",
851
                "answer": "<answer 5>"
852
            }}
853
854
855
    }}
856
    Only output the JSON. No additional commentary or explanations.
857
```

User prompt for generating CQA data:

```
859
860
    Generate 5 developmentally appropriate reading comprehension question-answer pairs
        based on the following input:
862
863
864
    - Text: {output}
    - Age Group: {age_group}
865
    - Stage: {stage}
866
867
868
    - Consider the developmental stage ({stage}) and age group ({age_group}) when
         crafting question complexity and answer expectations
870
    - Create questions that test different types of comprehension appropriate to the
871
         developmental level
872
873
    - **Ensure all questions can be answered directly from the provided text**
    - Generate authentic child responses that demonstrate comprehension at the target
874
875
876
    - Use vocabulary and sentence structures appropriate to the age group
877
    - Create a mix of question types that genuinely assess understanding of the text
878
879
    Output strictly in this format:
    "'json
880
    {{
881
882
        "question_answer_pairs": [
883
            {{
```

```
"question": "<question 1>",
884
                  "answer": "<answer 1>"
885
             }},
886
             }}
887
                 "question": "<question 2>",
888
                 "answer": "<answer 2>"
889
             }},
890
891
                 "question": "<question 3>",
892
893
                  "answer": "<answer 3>"
894
             }},
895
                 "question": "<question 4>",
896
                 "answer": "<answer 4>"
897
             }},
898
899
                  "question": "<question 5>",
900
                  "answer": "<answer 5>"
901
             }}
902
903
    }}
904
     ""
985
```

E.6 CSQA

System prompt for generating CSQA data:

```
908
909
    You are an AI model generating training data to help language models simulate human
910
911
        developmental skills at various stages from early childhood through early
         adolescence.
912
913
    Your task is to create 3 skill-based instruction-response pairs between an educator
914
915
         and a child that use a provided text as context to test specific developmental
916
         skills, rather than simple reading comprehension.
917
918
    Strictly follow these guidelines:
919
    1. **Developmental Appropriateness by Stage:**
920
       - Stage 0 (Age 5): Simple vocabulary, short sentences, concrete thinking, present
921
           -focused, immediate experiences
922
       - Stages 1-3 (Ages 6-8): Basic past/future concepts, simple reasoning, familiar
923
924
           contexts, beginning abstract thought
       - Stages 4-6 (Ages 9-11): Complex reasoning, abstract thinking, varied sentence
925
           structures, hypothetical scenarios
926
927
       - Stages 7-9 (Ages 12-14): Sophisticated vocabulary, multiple perspectives,
           advanced abstract reasoning, nuanced responses
928
929
    2. **Skill-Based Instruction Creation:**
930
       - **Use the provided text as context, not as the primary focus**
931
       - Create instructions that test the specific skill, subskill, goal, and indicator
932
933
       - Instructions should prompt the child to demonstrate the target skill using
934
           elements from the text
935
936
       - Avoid simple recall questions - focus on skill application, analysis, synthesis
937
           , or evaluation
       - Vary instruction starters - avoid overusing "Imagine..." or "Tell me about..."
938
       - Include necessary context within the instruction if recall is required
939
       - Use developmentally appropriate language and concepts for the target stage
940
       - Make instructions engaging and thought-provoking for the age group
941
942
    3. **Response Generation:**
943
       - Create authentic child responses that clearly demonstrate the target indicator
944
945
       - Use vocabulary, sentence complexity, and reasoning appropriate to the
           developmental stage
946
```

```
- Include natural speech patterns and expressions typical of the age group
947
        - Ensure responses show genuine skill application, not just text recall
        - Responses should be verifiable through either:
949
         * Information provided in the instruction or text
950
951
          * Common world knowledge appropriate for the child's developmental level
          * Typical personal experiences for that age group
952
953
        - Avoid arbitrary claims or purely imaginative details unless the skill
            explicitly encourages creativity
954
955
     4. **Context Integration:**
956
957
        - Use the provided text as a springboard for skill demonstration
        - Connect text elements to real-world applications of the skill
958
        - Encourage children to apply their skills to analyze, extend, or relate to the
959
            text content
960
        - Ensure the skill being tested is meaningfully connected to the text context
961
962
     5. **Content Guidelines:**
963
        - **Purely verbal exchanges** - no references to physical objects, gestures, or
964
            non-verbal actions
965
966
        - No formatting (bold, italics, markdown)
        - Instructions should feel natural and appropriate for educational settings
967
        - Responses should sound natural and spontaneous, not rehearsed
968
969
        - Include appropriate emotional expressions and personal connections when
970
        - Ensure logical consistency between instruction and response
971
972
        - Focus on the skill demonstration rather than text comprehension
973
     6. **Quality Standards:**
974
        - The exchange must demonstrate clear alignment with the skill, subskill, goal,
975
            and indicator
976
        - Each instruction must clearly target the specific developmental parameters
977
978
            provided
         Instructions should be distinct from each other, testing different aspects of
979
980
            the same skill
        - Both instruction and response should feel authentic to a real classroom or
981
982
            learning interaction
983
        - Responses must demonstrate clear mastery or development of the target skill
984
        - The text should serve as meaningful context, not just background information
        - Avoid overly abstract concepts for younger stages or overly simple concepts for
985
             older stages
986
        - Ensure developmental appropriateness in both challenge level and expectations
987
     7. **Output Format: ** Strictly return the output in the following JSON structure:
989
     ""json
990
     }}
991
992
         "skill_based_pairs": [
            {{
993
                "instruction": "<instruction 1>",
994
                "response": "<response 1>"
995
            }},
996
            {{
997
                 "instruction": "<instruction 2>",
998
                 "response": "<response 2>"
999
            }},
1000
1001
             }}
                "instruction": "<instruction 3>",
1002
                "response": "<response 3>"
1003
1004
            }}
        ]
1005
     }}
1006
1007
     Only output the JSON. No additional commentary or explanations.
1889
```

User prompt for generating CSQA data:

```
1011
1012
     Generate 3 developmentally appropriate skill-based instruction-response pairs based
1013
          on the following input:
1014
     - Text: {output}
1015
1016
       Age Group: {age_group}
     - Stage: {stage}
1017
     - Skill: {skill}
1018
1019
     - Sub-skill: {subskill}
     - Goal: {goal}
1020
     - Indicator: {indicator}
1021
1022
     Instructions:
1023
     - Consider the developmental stage ({stage}) and age group ({age_group}) when
1024
1025
          crafting instruction complexity and response expectations
     - Use the provided text as context to create instructions that test the specific
1026
          skill ({skill}) and subskill ({subskill})
1027
     - Create instructions that elicit demonstration of the goal ({goal}) and indicator
1028
          ({indicator})
1029
     - **Focus on skill application and demonstration, not text comprehension **
1030
     - Generate authentic child responses that show clear mastery of the target skill at
1031
          the developmental stage
1032
1033
      - Use vocabulary and sentence structures appropriate to the age group
1034
     - Create 3 distinct instructions that test different aspects of the same skill
1035
     Output strictly in this format:
1036
     ""json
1037
     {{
1038
1039
         "skill_based_pairs": [
1040
                 "instruction": "<instruction 1>",
1041
1042
                 "response": "<response 1>"
             }},
1043
             }}
1044
                 "instruction": "<instruction 2>",
1045
                 "response": "<response 2>"
1046
             }},
1047
1048
             {{
                 "instruction": "<instruction 3>",
1049
                 "response": "<response 3>"
1050
1051
             }}
1052
     }}
1053
     ""
1855
```

E.7 IR

1056

1057

System prompt for generating IR data:

```
1058
     You are an AI model generating training data to help language models simulate human
1059
         developmental skills at various stages from early childhood through early
1060
         adolescence.
1061
1062
1063
     Your task is to create realistic instruction-response pairs between an educator and
         a child, based on developmental indicators, skills, and a tuple of word and its
1064
          part of speech.
1065
1066
     Strictly follow these guidelines:
1067
1068
1069
     1. **Developmental Appropriateness by Stage:**
        - Stage 0 (Age 5): Simple vocabulary, short sentences, concrete thinking, present
1070
            -focused, immediate experiences
1071
1072
        - Stages 1-3 (Ages 6-8): Basic past/future concepts, simple reasoning, familiar
            contexts, beginning abstract thought
1073
```

```
1074
        - Stages 4-6 (Ages 9-11): Complex reasoning, abstract thinking, varied sentence
1075
            structures, hypothetical scenarios
        - Stages 7-9 (Ages 12-14): Sophisticated vocabulary, multiple perspectives,
1076
            advanced abstract reasoning, nuanced responses
1077
1078
     2. **Instruction Creation:**
1079
        - Use the provided word and its part of speech to meaningfully inspire the
1080
            interaction topic
1081
        - **Ensure the topic aligns with the Text Type Template (instruct_template)**
1082
1083
        - Craft prompts that naturally elicit demonstration of the specific indicator and
1084
             skill
        - Vary instruction starters - avoid overusing "Imagine..." or "Tell me about..."
1085
        - Include necessary context within the instruction if recall is required
1086
        - Use developmentally appropriate language and concepts for the target stage
1087
        - Make instructions engaging and thought-provoking for the age group
1088
1089
     3. **Response Generation:**
1090
        - Create authentic child responses that clearly demonstrate the target indicator
1091
        - Use vocabulary, sentence complexity, and reasoning appropriate to the
1092
1093
            developmental stage
        - Include natural speech patterns and expressions typical of the age group
1094
        - Ensure responses are verifiable through either:
1095
          * Information provided in the instruction
1096
          * Common world knowledge appropriate for the child's developmental level
1097
1098
          * Typical personal experiences for that age group
        - Avoid arbitrary claims or purely imaginative details unless storytelling is
1099
            explicitly encouraged
1100
1101
     4. **Content Guidelines:**
1102
        - **Purely verbal exchanges** - no references to physical objects, gestures, or
1103
            non-verbal actions
1104
        - No formatting (bold, italics, markdown)
1105
        - Responses should sound natural and spontaneous, not rehearsed
1106
        - Include appropriate emotional expressions and personal connections when
1107
1108
1109
        - Ensure logical consistency between instruction and response
1110
     5. **Quality Standards:**
1111
        - The exchange must demonstrate clear alignment with the skill, subskill, goal,
1112
            and indicator
1113
        - Both instruction and response should feel authentic to a real classroom or
1114
            learning interaction
1115
        - Avoid overly abstract concepts for younger stages or overly simple concepts for
1116
             older stages
1117
        - Ensure the selected word meaningfully influences the dialogue topic
1118
     6. **Output Format:** Strictly return the output in the following JSON structure:
1120
     ""json
1121
1122
     }}
         "instruction": "<instruction>",
1123
1124
         "response": "<response>"
1125
     }}
     ""
1126
     Only output the JSON. No additional commentary or explanations.
1128
```

User prompt for generating IR data:

```
Generate a developmentally appropriate instruction-response pair based on the
following input:

1132
1133
1134 - ID: {id}
1135 - Indicator: {indicator}
1136 - Skill: {skill}
1137 - Sub-skill: {subskill}
1138 - Goal: {goal}
```

```
- Age Group: {age_group}
1139
     - Stage: {stage}
1140
     - Text Type Template: {instruct_template}
1141
     - (Word, Part of speech): {word_list}
1142
1143
     Instructions:
1144
     - Consider the developmental stage ({stage}) and age group ({age_group}) when
1145
          crafting language complexity and content themes
1146
     - Use the selected word to meaningfully inspire the interaction topic **that fits
1147
1148
          the Text Type Template**
1149
     - Create an engaging instruction that naturally elicits demonstration of the target
1150
         indicator
     - Generate an authentic child response that clearly shows mastery of the skill and
1151
1152
     - Ensure the exchange feels natural and appropriate for a real educational
1153
          interaction
1154
1155
     Output strictly in this format:
1156
     ""json
1157
     {{
1158
         "instruction": "<instruction>",
1159
         "response": "<response>"
1160
     }}
1161
1163
```

E.8 Evaluating CQA

1164

1165

System prompt for evaluating trained model's response for questions from CQA:

```
You are a developmental expert evaluating how well a child's answer to a reading
1167
1168
          comprehension question reflects appropriate understanding and reasoning for a
          specific developmental stage.
1169
1170
1171
     You will receive:
     - The original **context** paragraph
1172
1173
     - A **question** based on the context
1174
      The child's **answer** to the question
     - The child's **developmental stage** (09)
1175
     - The child's **age group** (e.g., '05', '511', '1114')
1176
1177
1178
     Your job is to:
1179
     1. **Rate the childs answer on a scale from 1 to 5**, using the following criteria:
        - **5 Excellent:** Fully correct, precise, and well-formed for the stage. Shows
1180
            strong comprehension and reasoning.
1181
1182
        - **4 Strong:** Mostly correct and appropriate; may have minor phrasing issues
            or slight gaps in reasoning.
1183
        - **3 Adequate: ** Understands the gist but may be vague, partially incorrect, or
1184
             simplistic for the stage.
1185
        - **2 Limited: ** Misunderstands part of the question or context; reasoning is
1186
1187
            weak or off-track.
        - **1 Inadequate: ** Confused, incorrect, or clearly not appropriate for the
1188
1189
1190
1191
     2. **Consider developmental expectations** for language and reasoning:
        - **Stage 0 (Age 5):** Very basic phrases, literal recall, present-focused
1192
1193
            answers
        - **Stages 13 (Ages 68):** Simple reasoning, sequencing, basic cause-effect,
1194
            clear answers
1195
        - **Stages 46 (Ages 911): ** Logical inference, comparative language, clear
1196
            justification
1197
        - **Stages 79 (Ages 1214): ** Abstract reasoning, complex ideas, nuanced
1198
            explanations
1199
1200
    3. **Evaluate:**
1201
```

```
- Does the childs answer meaningfully address the question using the provided
1202
            context?
1203
        - Is the reasoning and language appropriate for the stage?
1204
        - Does it reflect comprehension of the text and question?
1205
1206
     4. **Output Format:**
1207
     Only return the following dictionary:
1208
     ""json
1209
     }}
1210
1211
         "rating": <integer from 1 to 5>,
         "explanation": "<23 sentence rationale>"
1212
     }}
1213
1214
     Do not add any other text or formatting. Only return the JSON object.
1215
```

1217 User prompt for evaluating trained model's response for questions from CQA:

```
1218
     Evaluate the childs answer to a reading comprehension question. Consider the context
1219
1220
           and the developmental stage.
1221
1222
     Context:
     {context}
1223
1224
1225
     Question:
     {question}
1226
1227
1228
     Answer:
     {answer}
1229
1230
     Stage: {stage}
1231
     Age group: {age_group}
1232
     Index: {q_index}
1233
     **Output Format:**
1235
     "''json
1236
1237
     {{
1238
          "rating": <integer from 1 to 5>,
          "explanation": "<23 sentence rationale>"
1239
     }}
1240
     ""
1342
```

E.9 Evaluating CSQA

1243

1244

System prompt for evaluating trained model's response for questions from CSQA:

```
1245
     You are a developmental expert evaluating how well a child's response demonstrates a
1246
           specific developmental skill at a given stage, using a provided instruction
1247
          and background text.
1248
1249
     You will receive:
1250
     - A short **text** (used as context for the instruction)
1251
     - A **skill-based instruction** given to the child
1252
1253
     - The childs **response**
     - The childs **developmental stage** (09)
1254
     - The childs **age group** (e.g., '05', '511', '1114')
1255
     - The **target skill**, **subskill**, **goal**, and **indicator** that the
1256
1257
          instruction was designed to assess
1258
1259
     Your job is to:
     1. **Rate the child's response on a scale from 1 to 5**, using these criteria:
1260
        - **5 Excellent:** Fully demonstrates the targeted skill/indicator with clarity
1261
1262
            and developmental appropriateness. Strong reasoning, appropriate expression,
1263
            and alignment with instruction.
```

```
- **4 Strong: ** Mostly appropriate and well-formed. Some minor gaps in
1264
            completeness, precision, or phrasing, but shows the intended skill.
1265
        - **3 Adequate: ** Response attempts the skill but may be vague, simplistic, or
1266
            only partially aligned with the goal/indicator.
1267
        - **2 Limited: ** Weak or unclear demonstration of the skill. Response is
1268
            partially off-track, underdeveloped, or barely relevant.
1269
1270
        - **1 Inadequate: ** Fails to demonstrate the intended skill. Response is
            irrelevant, confusing, or clearly inappropriate for the stage.
1271
1272
1273
     2. **Use stage-specific developmental expectations**:
1274
        - **Stage 0 (Age 5):** Short, concrete, present-focused responses with simple
1275
            vocabulary
        - **Stages 13 (Ages 68): ** Clear expression of ideas, simple cause-effect,
1276
1277
            emotional awareness, basic reasoning
        - **Stages 46 (Ages 911):** Logical structure, hypothetical thinking, connections
1278
             to personal experience, comparisons
1279
        - **Stages 79 (Ages 1214): ** Advanced abstraction, multiple perspectives,
1280
            justification, nuanced expression
1281
1282
1283
     3. **Evaluate:**
        - Does the childs response meaningfully follow the instruction?
1284
        - Does it demonstrate the **targeted skill and indicator**?
1285
1286
        - Is the language, reasoning, and expression developmentally appropriate for the
1287
        - Is the response authentic and logically consistent with the instruction and the
1288
1289
             context text?
1290
     4. **Output Format:**
1291
1292
     Return only the following dictionary:
     "json
1293
     {{
1294
         "rating": <integer from 1 to 5>,
1295
         "explanation": "<23 sentence rationale>"
1296
     }}
1297
1298
     Do not add any other text or formatting. Only return the JSON object.
1388
```

User prompt for evaluating trained model's response for questions from CSQA:

```
Evaluate the child's response to a skill-based instruction using the provided text
1303
          and developmental context. Focus on how well the response demonstrates the
1304
          intended skill.
1305
1306
     Context:
1307
     {context}
1308
1309
     Instruction:
1310
     {instruction}
1311
1312
     Response:
1313
     {response}
1314
1315
     Stage: {stage}
1316
     Age group: {age_group}
1317
1318
     Skill: {skill}
     Subskill: {subskill}
1319
     Goal: {goal}
1320
1321
     Indicator: {indicator}
1322
     Index: {q_index}
1324
     Output format:
     "; json
1325
     {{
1326
1327
          "rating": <integer from 1 to 5>,
          "explanation": "<23 sentence rationale>"
1328
```

```
1329 | }}
'''
```

E.10 Evaluating IR

1332

1384

System prompt for evaluating trained model's response for questions from IR:

```
1334
     You are a developmental expert rating how well a child's response to a prompt
1335
1336
          demonstrates age-appropriate reasoning and language for a given developmental
1337
          stage.
1338
1339
     You will receive:
     - An **instruction** given to the child
1340
1341
     - The child's **response**
     - The child's **developmental stage** (09)
1342
     - The child's **age group** (e.g., '05', '511', '1114')
1343
1344
1345
     Your job is to:
     1. **Rate the response on a scale from 1 to 5**, using the following criteria:
1346
        - **5 Excellent: ** The response fully addresses the instruction with clear,
1347
1348
            developmentally appropriate reasoning and language. It meets expectations for
             the stage with no major issues.
1349
        - **4 Strong: ** Mostly appropriate and coherent; minor gaps in clarity, depth,
1350
            or completeness.
1351
        - **3 Adequate: ** A reasonable attempt that partially addresses the instruction;
1352
             may be vague, brief, or contain small misunderstandings.
1353
1354
         **2 Limited:** Weak or underdeveloped response; minimal reasoning or limited
            relevance to the instruction.
1355
        - **1 Inadequate: ** Response is off-topic, confusing, or clearly inappropriate
1356
1357
            for the stage.
1358
     2. **Use stage-specific developmental expectations**:
1359
        - **Stage 0 (Age 5):** Very simple sentences, concrete ideas, focused on here and
1360
1361
             now
        - **Stages 13 (Ages 68): ** Simple reasoning, some past/future thinking, familiar
1362
1363
            examples
        - **Stages 46 (Ages 911): ** Logical structure, comparisons, abstract or
1364
            hypothetical reasoning
1365
        - **Stages 79 (Ages 1214): ** Nuanced reasoning, multi-step thinking, advanced
1366
            vocabulary
1367
1368
     3. **Evaluate:**
1369
        - Does the childs response meaningfully address the instruction?
1370
        - Is the language and reasoning developmentally appropriate for the stage?
1371
        - Is the response authentic and logically consistent?
1372
1373
     4. **Output Format:**
1374
     Only return the following dictionary:
1375
     ";json
1376
     {{
1377
1378
         "rating": <integer from 1 to 5>,
         "explanation": "<23 sentence rationale>"
1379
1380
     }}
1381
     Do not add any other text or formatting. Only return the JSON object.
1383
```

User prompt for evaluating trained model's response for questions from IR:

```
Evaluate the child's response to the instruction below based on the developmental
stage and age group. Return a numerical rating (15) and a short explanation.

Instruction: {instruction}
Response: {response}
Stage: {stage}
```

```
Age group: {age_group}
Index: {q_index}
1392
1393
1394
       **Output Format:**
Only return the following dictionary:
'''json
1395
1396
1397
       {{
1398
             "rating": <integer from 1 to 5>,
"explanation": "<23 sentence rationale>"
1399
1400
       }}
1401
1483
```