

From Amateur to Master: Infusing Domain Knowledge into LLMs via Automated Curriculum Learning

Anonymous ACL submission

Abstract

Large Language Models (LLMs) excel at general tasks but underperform in specialized domains like economics and psychology, which require deep, principled understanding. To address this, we introduce ACER (Automated Curriculum-Enhanced Regimen) a framework for targeted domain knowledge infusion that combines structured synthetic corpus generation with curriculum-aligned continual pre-training. ACER synthesizes textbook-style curriculum with complementary question-answer pairs guided by Bloom’s taxonomy, enabling systematic coverage and progressive cognitive difficulty. The resulting synthetic corpus is used to drive curriculum-aligned continual pre-training, rather than relying on unstructured or naively mixed data. Experiments on Llama 3.2 (3B and 1B) show consistent improvements on five challenging MMLU subdomains, with gains of up to 5 percentage points in particularly difficult areas such as microeconomics and a macro-average improvement of about 3 points across target domains. Importantly, ACER preserves performance on non-target domains and often yields modest positive transfer. Beyond MMLU, ACER improves performance on knowledge-intensive benchmarks such as ARC and GPQA by over 2 absolute points, while maintaining stable performance on general reasoning tasks. Overall, ACER provides a scalable approach for infusing principled domain expertise into general-purpose LLMs without sacrificing their breadth.

1 Introduction

Large Language Models (LLMs) have achieved remarkable success across a wide range of natural language processing (NLP) tasks, including open-domain question answering, summarization, reasoning, and code generation, largely driven by scaling both model size and training data (Kaplan et al., 2020). However, this broad competency masks a critical weakness. While state-of-the-art models

excel in general tasks, they falter in specialized domains that demand deep, principled understanding (He et al., 2025; Zhang et al., 2025b). This gap is highlighted in knowledge benchmarks like MMLU (Hendrycks et al., 2021), where performance in niche sub-domains like virology degrades, compared to broader medicine, revealing a gap between general knowledge and deep expertise.

This performance gap stems from the nature of LLM pretraining corpora, which are dominated by general web text and underrepresent specialized knowledge (Najem-Meyer et al., 2025). Even when domain-specific data is present, it often lacks the methodical exposition and progressive knowledge scaffolding of expert materials like textbooks or lecture notes. Consequently, even large-scale models struggle with the technical terminology and hierarchical concepts essential for expert-level performance (Mai et al., 2024). Several strategies have been explored to address this gap. Instruction tuning improves alignment but not knowledge depth (Wei et al., 2022), while domain-specific pretraining often suffers from data scarcity and degradation of general capabilities (Gururangan et al., 2020; Béthune et al., 2025). Synthetic data generation has recently emerged as a promising alternative, with both instruction-centric methods (e.g., Self-Instruct (Wang et al., 2023), GLAN (Li et al., 2025)) and large-scale synthetic pretraining efforts (e.g., Phi (Abdin et al., 2024), Cosmopedia (Ben Allal et al., 2024)) demonstrating clear benefits. However, these methods often lack systematic coverage of domain concepts or structured alignment, limiting their effectiveness in building principled domain expertise. This leaves a critical need for an approach that can instill deep, methodical knowledge without undermining the broad capabilities that make LLMs so powerful.

To address this challenge, we introduce ACER (Automated Curriculum-Enhanced Regimen), a framework designed to infuse LLMs with domain

085 expertise without sacrificing their broad applicabil- 137
086 ity. Our methodology consists of two core com- 138
087 ponents: a process for systematically generating 139
088 expert “study materials” and a novel training re- 140
089 gimen for the model to absorb them. It begins 141
090 by generating a detailed table of contents (ToC) 142
091 that serves as the blueprint for textbook construc- 143
092 tion. The content is then generated section by 144
093 section, resulting in a comprehensive synthetic 145
094 textbook. Guided by Bloom’s taxonomy (Bloom 146
095 et al., 1956), ACER then complements synthetic 147
096 textbooks with exam-style QA pairs to ensure sys- 148
097 tematic coverage and progressive difficulty. To 149
098 reflect educational progression, we generate four 150
099 versions of each textbook tailored to different au- 151
100 diences: high school, undergraduate, graduate, 152
101 and researcher. The resulting synthetic corpora 153
102 combine topical breadth with structured progres- 154
103 sion, enabling LLMs to acquire principled domain- 155
104 specific knowledge through systematic curriculum 156
105 progression that parallels human educational devel-
106 opment. This synthetic curriculum is then used to
107 continually pretrain a foundational LLM. We em-
108 ploy curriculum-aligned training schedules across
109 cognitive and content dimensions that strategically
110 mix the new expert corpora with general-domain
111 data, enabling the model to gain deep expertise
112 while retaining its broad capabilities.

113 Our evaluation methodology begins by systemat-
114 ically identifying a model’s most significant knowl-
115 edge gaps. To this end, we benchmarked the
116 Llama 3.2 3B and 1B (Grattafiori et al., 2024) "stu-
117 dent" models against their Llama 3.1 8B "teacher"
118 across all 56 MMLU domains. The five domains
119 exhibiting the most severe performance degrada-
120 tion were then selected as the proving ground for
121 ACER’s ability to build targeted expertise. We gen-
122 erated synthetic book corpora for these domains
123 and continually pretrained the baseline models on
124 this corpus, mixing synthetic data with general-
125 domain replay data. To ablate different curriculum
126 effects, we experimented with multiple scheduling
127 strategies, including cognitive ordering (textbook
128 → easy QA → hard QA) and persona-based con-
129 tent ordering (high school → undergraduate →
130 graduate → researcher).

131 **Strong Results:** ACER consistently outperformed
132 pretrained baselines, with cognitive and content
133 based curriculum yielding macro-average gains of
134 about 3 percentage points in the target domains.
135 Specifically in challenging areas such as microe-
136 conomics, ACER improved accuracy by nearly 5

137 points. Additionally, performance on non-target 138
139 domains is preserved relative to the pretrained base- 140
141 line, with small positive changes observed in sev- 142
143 eral cases. Beyond MMLU, we further evaluated 144
145 ACER on widely used benchmarks such as ARC 146
147 (Clark et al., 2018), GPQA (Rein et al., 2023), 148
149 AGIEval (Zhong et al., 2024), GSM8K (Cobbe 150
151 et al., 2021), and HellaSwag (Zellers et al., 2019). 152
153 The knowledge-infused ACER-trained models im- 154
155 proved by more than 2 absolute points in ARC 156
157 and GPQA, both of which emphasize knowledge 158
159 recall and domain understanding, while maintain- 160
161 ing stable performance on general reasoning, arith- 162
163 metic, and common sense tasks such as AGIEval, 164
165 GSM8K, and HellaSwag. These results demon- 166
167 strate that continual pretraining with ACER not only 168
169 enhances specialized knowledge, but also preserves 170
171 broad capabilities, providing a scalable recipe for 172
173 closing domain gaps in LLMs. 174
175

176 Our contributions are:

- 177 1. **Curriculum-Aligned Knowledge Infusion** 178
179 **Framework:** We introduce ACER, a frame- 180
181 work for targeted domain knowledge infusion 182
183 that integrates structured synthetic textbook- 184
185 style curriculum with complementary exam- 186
187 style question-answer pairs across multiple 188
189 educational levels. Unlike prior synthetic pre- 189
190 training approaches that rely on broad or un- 190
191 structured data, ACER aligns content coverage 191
192 and cognitive progression to support princi- 192
193 pled domain learning in LLMs (Section 3). 193
- 194 2. **Effective Curriculum Learning Regimen:** 194
195 We design and evaluate curriculum learning 195
196 strategies (Section 3.1) that control both con- 196
197 tent ordering and cognitive difficulty during 197
198 continual pretraining. With these curriculum 198
199 designs, ACER achieves consistent improve- 199
200 ments over pretrained Llama 3.2 baselines, 200
201 with a macro-average gain of about 3 points 201
202 across target MMLU domains and particularly 202
203 strong improvements of up to 5 points in mi- 203
204 croeconomics (Section 4.2). 204
- 205 3. **Robust Generalization:** We demonstrate that 205
206 ACER generalizes beyond in-domain tasks in 206
207 MMLU, yielding over 2 absolute point im- 207
208 provements on knowledge-intensive bench- 208
209 marks such as ARC and GPQA, while pre- 209
210 serving capabilities on general reasoning, 210
211 arithmetic, and commonsense tasks such as 211
212 AGIEval, GSM8K, and HellaSwag. These 212
213

187 results indicate that structured, curriculum- 237
188 based continual pretraining enables domain 238
189 specialization without degrading general- 239
190 purpose performance (Section 4.3). 240

191 2 Related Work 243

192 LLMs have been very useful in general tasks, but 244
193 they lag considerably in niche domains that de- 245
194 mand deep, principled understanding (He et al., 246
195 2025; Zhang et al., 2025b). This performance gap 247
196 is consistently reflected in benchmarks such as 248
197 MMLU (Hendrycks et al., 2021), which span a 249
198 wide range of specialized domains. One contribut- 250
199 ing factor is that specialized domains remain un- 251
200 derrepresented in the pretraining corpora (Najem- 252
201 Meyer et al., 2025). Domain-adaptive pretraining 253
202 has been effective in improving LLMs in such set- 254
203 tings. Don't Stop Pretraining (Gururangan et al., 255
204 2020) demonstrated that continual pretraining on 256
205 domain-specific corpora can significantly improve 257
206 downstream task performance. (Kerner, 2024) fur- 258
207 ther showed that even compact, specialized mod- 259
208 els can achieve competitive accuracy in-domain 260
209 compared to much larger general-purpose LLMs. 261
210 One reason, as suggested by (Mai et al., 2024), is 262
211 that large general-purpose models struggle with 263
212 domain-specific reasoning and hierarchical con- 264
213 cepts essential for expert performance. More re- 265
214 cent efforts, such as PreparedLLM (Chen et al., 266
215 2024), combine instruction-based pretraining with 267
216 domain adaptation under a structured, curriculum- 268
217 style regime to further improve specialization. De- 269
218 spite these advances, domain-adaptive pretraining 270
219 faces fundamental challenges. High-quality cu- 271
220 rated corpora are often scarce or subject to licens- 272
221 ing restrictions, making it difficult to scale adapta- 273
222 tion to multiple domains (Wu et al., 2025). More- 274
223 over, continual pretraining on narrow domains can 275
224 lead to degradation of general capabilities, as high- 276
225 lighted by recent studies (Béthune et al., 2025; 277
226 Huang et al., 2024). Thus, while domain-adaptive 278
227 pretraining improves specialization, it remains con-
228 strained by data scarcity and forgetting risks, moti-
229 vating our work.

230 Synthetic data generation has recently emerged
231 as a promising alternative with large-scale cor-
232 pora built from scratch. For instance, Cosmo-
233 pedia (Ben Allal et al., 2024) uses carefully de-
234 signed multi-stage prompts to generate diverse
235 open textbook-style pretraining corpora, while
236 the Phi-4 models (Abdin et al., 2024) leverage

multi-agent, multi-stage prompting pipelines to pro-
duce synthetic datasets spanning hundreds of bil-
lions of tokens across diverse domains and skills.
These large-scale efforts complement earlier ap-
proaches such as Self-Instruct (Wang et al., 2023)
and GLAN (Li et al., 2025), which demonstrate that
high-quality synthetic instruction-response pairs
can improve model generalization, providing scal-
able alternatives to human-annotated corpora for
post-training. Beyond general-purpose synthesis,
some methods have lately focused on targeted do-
main adaptation. For instance, (Arannil et al.,
2024) mine domain-related subsets from large
web datasets, while (Yang et al., 2024) generate
synthetic text by extracting salient domain enti-
ties from documents and constructing connections
among them. While such approaches highlight the
scalability and diversity benefits of synthetic data,
they typically lack systematic coverage of domain
concepts or curriculum-aligned progression, which
are crucial for instilling principled and progressive
domain expertise in LLMs.

The order in which training data is presented
has a strong influence on model performance, as
first demonstrated by (Bengio et al., 2009). Re-
cent work has extended this idea to LLMs. (Zhang
et al., 2025a) study multiple curriculum strate-
gies for pretraining, including vanilla ordering,
pacing-based sampling, and interleaving, show-
ing that thoughtful sequencing can improve train-
ing efficiency and downstream accuracy. Comple-
mentarily, (Lee et al., 2024) introduce curriculum
instruction tuning, where sequencing instruction-
response pairs by difficulty improves results across
diverse benchmarks without additional computa-
tional costs. However, such strategies remain
largely underexplored in the context of domain
knowledge infusion or synthetic textbook-style cor-
pora. In this work, we investigate curriculum
scheduling as a principal mechanism to instill pro-
gressive knowledge in domain-adaptive pretrain-
ing.

279 3 Automated Curriculum-Enhanced 280 Regimen - ACER

281 ACER is a multi-stage pipeline designed to trans-
282 form high-level learning goals into a structured syn-
283 thetic training corpus. As shown in Figure 1, the
284 process begins by capturing domain intent and au-
285 dience context, expands this information into a de-
286 tailed outline, and then produces textbook-style sec-

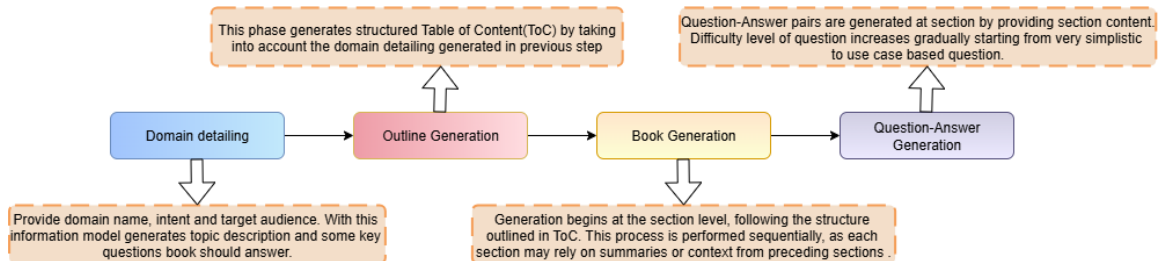


Figure 1: ACER - pipeline for synthetic book corpus generation

tions followed by section-aware question–answer pairs. Each stage feeds the next, encouraging consistency, progressive depth, and factual grounding.

Domain Detailing: The first stage of ACER is domain detailing. The goal of this stage is to establish a structured foundation for synthetic corpus generation by capturing domain-specific knowledge requirements, audience context, and intent. This stage aims to formalize the scope, granularity, and emphasis of the target content before automated generation begins, ensuring that all downstream artifacts (textbooks, question–answer pairs, and curricula) are aligned with both educational principles and application needs.

The process begins with three primary inputs: **Domain or Topic Name:** A concise label identifying the subject area, e.g., Anatomy, Microeconomics; **Intent:** A statement describing the purpose of the synthetic corpus. For instance, training domain experts, creating introductory learning materials, supporting professional certification preparation, etc.; and **Target Audience Metadata:** Learner profile specifying prior knowledge level (e.g., high school, undergraduate, graduate, or researcher) and contextual constraints (e.g., technical rigor, professional applications). To translate these inputs into a machine-usable representation, we employ a prompt-driven planning step in which a language model, acting as a “domain author”, generates three key artifacts:

- (1) **Domain Description:** A concise, yet comprehensive overview of the topic tailored to the target audience.
- (2) **Core Subtopics:** A list of areas essential for building expertise, ensuring systematic coverage
- (3) **Key Questions:** A set of 6-8 relevant questions that the textbook should answer, aligning with Bloom’s taxonomy objectives, such as comprehension, application, and synthesis.

The output of this stage is a JSON-encoded schema that captures the domain’s description,

subtopics, and key questions. This structured representation is editable, enabling subject-matter experts to iteratively refine content priorities before subsequent stages. Additional details about this stage are described in Appendix A.1. The prompt used to generate domain detailing are present in Appendix G (Code Block: 3)

Outline Generation: Once the domain schema is finalized in the domain detailing stage, the next step is outline generation, where the framework transforms structured topic metadata into a detailed and hierarchical Table of Contents (ToC). This step provides a precise blueprint for the creation of synthetic textbooks, ensuring both thematic coverage and logical progression of ideas. In addition to topic name, intent and target audience (as described in domain detailing), the outline generation process uses the following inputs: **Genre and Style Parameters:** content preferences, including tone, narrative voice, and language style, to ensure readability and audience alignment; and **Domain Schema:** description, core subtopics, and key questions generated in the domain detailing stage. Collectively, these attributes act as guidance signals for the language model, shaping the structure and content depth of the generated outline. More details about the outline generation phase is described in Appendix A.2. The prompt used to generate the ToC can be found in the Appendix G (Code Block: 4)

Synthetic Content Generation: Following the generation of the Table of Contents (ToC), the next step involves the generation of synthetic textbook-style content. The ToC is used not only as a hierarchical list, but as a structured blueprint that is systematically traversed and expanded into full-fledged instructional material. This process operates at the *section-level granularity*, ensuring that the generated text is pedagogically cohesive while remaining contextually aligned with the surrounding chapters and sections. A detailed description

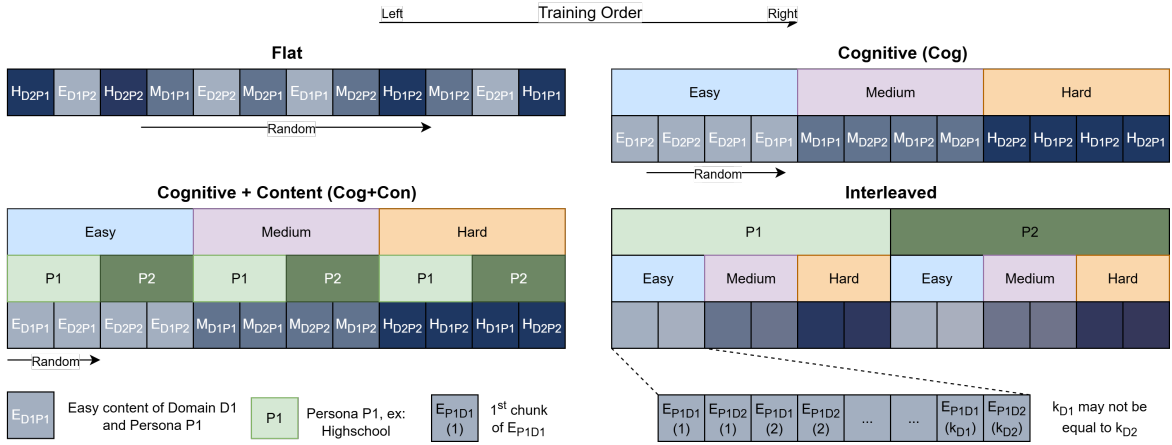


Figure 2: Different training schedules. Personas $P1$ and $P2$ are shown here only for illustration and do not represent the full set of personas. Persona $P2$ lies higher than $P1$ on the cognitive axis (e.g., $P2$ =Undergraduate, $P1$ =high school). $E_{P1D1}(1)$ denotes the first chunk of E_{P1D1} , which is formed by combining multiple consecutive sections.

of the procedure and structure is described in Appendix A.3.

Following the synthesis of section-level content, we extend the generation pipeline to incorporate question–answer (QA) pairs derived from the generated text. The rationale for this stage is rooted in pretraining needs: question–answer pairs have been shown to be particularly effective in enhancing reasoning capabilities and comprehension during large-scale model training (Cheng et al., 2024). We generate QA pairs with two sequential stages:

- Question Generation:** The data generation model is framed as a subject-matter expert and tasked with generating self-contained, educational questions tied to the section content. For the *first question*, the prompt emphasizes simplicity, typically asking for a factual recall or definition-based understanding. For subsequent questions, the difficulty gradually escalates, following a cognitive progression from comprehension to interpretation, analysis, and application (Bloom et al., 1956). Furthermore, the prompt enforces strict relevance and audience alignment, while prohibiting extraneous labels or formatting.
- Answer Generation:** Once the question is obtained, a separate model call generates a detailed, educational, and self-contained answer. The answer is grounded strictly in the section content, ensuring factual alignment and eliminating hallucinations. The model is prompted to produce responses that are rich in detail, employ clear language, and illustrate concepts with examples. The tone is maintained at an

educational and explanatory level, suitable for the designated target audience.

The generated question answer pairs are then divided into easy and hard subsets according to their difficulty. The different prompts used for generating section content, question and answer pairs can be found in Appendix G (Code Blocks: 5, 6, 7 respectively).

3.1 ACER - Curriculum Scheduling

We experiment with four curriculum schedules for incorporating synthetic book corpus into continual pretraining (refer Figure 2):

- Flat:** All data (books, easy QA, hard QA) from all domains and personas are presented together without any ordering. The data can come in any order, irrespective of their difficulty level.
- Cognitive (Cog):** Data is structured in increasing cognitive difficulty: Books \rightarrow Easy QA \rightarrow Hard QA. Inside each of these sections, there is no restriction on the ordering of the domain or persona. Corpora from different domains and personas are mixed randomly.
- Cognitive + Content (Cog+Con):** Extends the cognitive schedule by introducing content-based ordering: High school \rightarrow Undergraduate \rightarrow Graduate \rightarrow Researcher. Again, domains can appear in any order and do not need to follow the same sequence across different personas.

4. **Interleaved:** Inspired by Lee et al. (2024), the data is interleaved across domains at the chapter–section level (e.g., Chapter 1, Section 1 from Domain 1 → Chapter 1, Section 1 from Domain 2 → ...). This interleaving is applied separately for books, easy QA, and hard QA, while preserving the persona order within each category. The model first encounters high-school content, followed by undergraduate, graduate, and researcher material, with section-level interleaving for each persona. This arrangement prevents the model from seeing all data from a single domain consecutively. Instead, it alternates between sections from multiple domains while respecting both cognitive difficulty and persona order, following a fixed cyclic pattern. This design enables us to test whether continual pretraining with structured sequencing of synthetic knowledge improves model quality.

4 Experiments

We experiment with Llama 3.2 3B and 1B, which are compact student models trained using supervision from larger Llama 3.1 teachers (8B and 70B).

To identify domains where domain-specific knowledge infusion is needed the most, we measured the per-domain accuracy gap between the 3B student and its 8B teacher on the full set of 56 MMLU tasks under 0-shot evaluation. Figure 3 shows the difference in accuracy between the student and teacher, with domains ranked by gap size. The largest regressions occurred in microeconomics, statistics, econometrics, mathematics, and psychology. These domains serve as our primary infusion targets. For consistency, we use this same set of domains for both the 3B and 1B students, while treating all other MMLU domains as non-targets.

Our synthetic book corpus consists of detailed, domain-specific books that span across diverse topics, each accompanied by exam-style question–answer (QA) pairs. The QA sets are divided into two difficulty levels: easy and hard. For each target domain, we generate synthetic corpora across four audience (persona) levels: high school, undergraduate, graduate, and researcher. To avoid benchmark contamination, we decontaminate the corpus by removing any text with cosine similarity greater than 0.9 to MMLU content; Appendix H details the decontamination process. Appendix B sum-

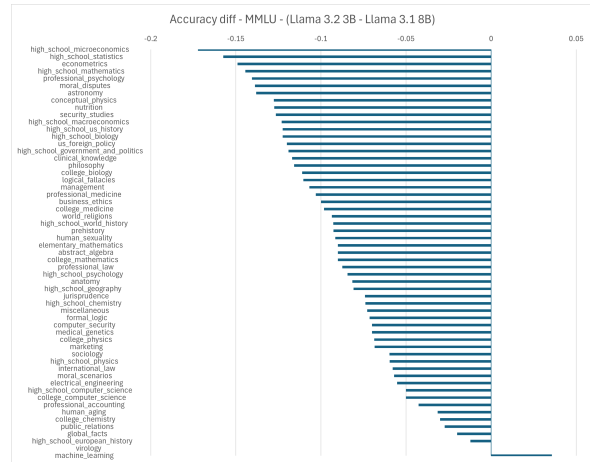


Figure 3: Target domains for knowledge infusion: identified using accuracy difference between the Llama 3.2 3B student model and its teacher in Llama 3.1 8B

marizes the token distribution across domains and audience levels: on average, ACER generates 6.4M tokens per domain, totaling about 32M tokens overall. Appendix E lists the generated artifacts for microeconomics domain.

4.1 Training Details

We initialize from the pretrained Llama 3.2 3B and 1B checkpoints and continually pretrain them on ACER generated synthetic book corpora using the standard next token prediction objective. To mitigate degradation of general capabilities, the synthetic in-domain corpus is mixed 1:1 with general replay data¹. The key hyper-parameters for training are set as: batch size = 512, maximum sequence length = 8192, learning rate = 2×10^{-5} with cosine decay to 2×10^{-6} , and a warm-up of 1% to peak learning rate. The replay data is drawn from pile knowledge, cosmopedia, ultratextbooks, and subsets of proof-pile-2 (openwebmath, algebraic stack, and arXiv). Using the curriculum schedules described in Section 3.1 (Flat, Cog, Cog+Con, and Interleaved), we evaluate the impact of data ordering on domain knowledge infusion.

For evaluation, we focus on five target MMLU subsets: high school microeconomics (MEco_{hs}), high school statistics (Stats_{hs}), econometrics (Econ), high school mathematics (Maths_{hs}), and professional psychology (Psych_p). We report both per-domain accuracy and the macro-average across these five subsets, denoted Macro_t. To track generalization, we also report macro-average accuracy

¹Appendix D.4 details our ablation study showing 1:1 ratio provides the right tradeoff

Model (Llama 3.2 3B)	MEco _{h,s}	Stats _{h,s}	Econ	Maths _{h,s}	Psych _p	Macro _t	Macro _{nt}
Pretrained baseline	0.5378	0.3796	0.3158	0.2704	0.5507	0.4108	0.5754
Flat _{gemini-2.0-flash}	0.5588	0.3843	0.3596	0.3111	0.5670	0.4362	0.5824
Cog _{gemini-2.0-flash}	0.5966	0.3796	0.3596	0.2963	0.5686	0.4401	0.5809
Cog+Con _{gemini-2.0-flash}	0.5840	0.4028	0.3509	0.2889	0.5768	0.4407	0.5821
Interleaved _{gemini-2.0-flash}	0.5798	0.3611	0.3509	0.2593	0.5605	0.4223	0.5766
Robustness to Synthetic Data Generators							
Cog+Con _{Phi-4}	0.5588	0.3704	0.3246	0.2963	0.5490	0.4198	0.5767
Cog+Con _{Llama-3.1-70B}	0.563	0.3565	0.307	0.2963	0.5408	0.4127	0.5749

Table 1: **Top Block:** Comparison of Llama 3.2 3B with ACER-trained variants under different curriculum schedules using synthetic data generated by Gemini 2.0 Flash. Domain-targeted synthetic textbooks with exam-style QA improve the pretrained baseline across curriculum schedules, with **Cog+Con** providing an improvement of **3.0** percentage points on target domains. Performance on non-target domains is preserved relative to the pretrained baseline, with small positive changes observed. **Bottom Block:** Robustness of the best-performing curriculum (Cog+Con) to different synthetic data generators. While generator choice modulates the magnitude of gains, curriculum-aligned training improves or preserves performance without degradation.

across the remaining 51 MMLU subsets, denoted as Macro_{nt}. Beyond MMLU, we evaluate the models on ARC, GPQA, AGIEval, GSM8K, and HellaSwag to assess broader reasoning and knowledge capabilities.

Positioning Relative to Synthetic-Pretraining Models:

Before presenting ACER results, we situate our study relative to recent pretraining efforts that expand knowledge coverage via large-scale synthetic data. We evaluate representative 1B-3B parameter models trained predominantly on broad synthetic corpora, specifically Phi 1.5, Phi 2, and Cosmopedia, on the same five MMLU subsets considered in this work. As detailed in Appendix C, these models exhibit low performance in specialized domains such as econometrics and professional psychology, indicating limited domain-specific knowledge despite their scale and diversity. Although these results are not directly comparable to ACER due to differences in architecture and training objectives, they suggest that general-purpose synthetic pretraining alone does not reliably close fine-grained domain gaps, motivating the targeted, curriculum-aligned knowledge infusion explored in ACER.

4.2 Results: with Llama 3.2 baselines

We generate synthetic book corpora for the identified target domains and continually pretrained the baseline Llama 3.2 3B model using different curriculum schedules. Table 1 compares these trained models with the pretrained 3B baseline. Across all schedules, ACER consistently improves perfor-

mance on the target domains. The **Flat** regime, which uniformly mixes books and QA pairs without ordering, yields a strong gain of +2.5 points in Macro_t over the baseline. Building on this, the **Cognitive (Cog)** schedule, which progresses from books → easy QA → hard QA, provides further incremental improvements. Extending this with persona-based content progression in **Cog+Con** delivers the best overall performance, indicating that curriculum design amplifies the benefits of synthetic corpora in continual pretraining.

In contrast, the **Interleaved** schedule, which mixes book sections and QA pairs across domains within training stages, underperforms relative to the other curriculum schedules, with particularly large drops in mathematics and statistics. This contrasts with Lee et al. (2024), where interleaving educational content improved performance of the model in the fine-tuning phase. We conjecture that, while interleaving may be effective for instruction-following fine-tuning, its fragmented sequencing dilutes the supervision signal in continual pretraining. For smaller-scale domain infusion tasks, such as ours, frequent task-switching introduced by interleaving likely overwhelms model capacity, leading to degraded performance rather than gains.

As shown in Table 1, domain knowledge infusion through targeted continual pretraining provides the largest benefits in microeconomics and econometrics. These niche areas are relatively underrepresented in the pretraining corpus of the 3B baseline, leaving substantial room for improvement. By contrast, domains such as statistics and mathematics are more prevalent in web-scale text, so

Benchmark	#shots	Setting	Llama 3.2 3B		Llama 3.2 1B	
			Pre-trained	Ours	Pre-trained	Ours
ARC Challenge	25	Acc (Weighted)	0.4701	0.4837	0.3652	0.3831
GPQA	0	Acc (Mean)	0.2656	0.2879	0.2545	0.2478
MMLU	5	Acc (Weighted)	0.5408	0.569	0.3099	0.3277
AGIEval	0	Acc (Weighted)	0.2255	0.2253	0.1867	0.1877
GSM8K	5	Strict (EM)	0.2805	0.2798	0.0667	0.0591
GSM8K	5	Flexible (EM)	0.2767	0.276	0.0644	0.0523
HellaSwag	0	Acc (Mean)	0.5522	0.5489	0.4777	0.4782

Table 2: Evaluation across benchmarks comparing our knowledge-infused Llama 3.2 models (3B and 1B) with their pretrained versions. The knowledge-infused models improve tasks requiring knowledge recall and understanding (ARC, GPQA, and MMLU) by more than 2 absolute percentage points, without regressing on general capabilities such as language understanding, reasoning, and mathematics (AGIEval, GSM8K, and HellaSwag).

the baseline already demonstrates strong competence, resulting in only modest gains. Overall, we observe an improvement of approximately 3 percentage points in Macro_t . Notably, our training regimen preserves performance on non-target domains (Macro_{nt}), with small positive changes observed in several cases. Appendix F shows an example where ACER enabled the model to correctly answer a microeconomics question that the pretrained baseline answered incorrectly.

We also evaluate the best-performing Cog+Con curriculum on Llama 3.2 1B, observing improvements of up to 2.7 points on target domains, with detailed results reported in Appendix D.3.

We further examine the sensitivity of ACER to the choice of synthetic data generator by re-generating the curriculum using different models, while fixing the curriculum schedule to the best-performing variant. As shown in the bottom block of Table 1, generator choice modulates the magnitude of gains, but curriculum-aligned training consistently improves or preserves performance relative to the pretrained baseline. A detailed analysis is provided in Appendix D.2.

4.3 Impact Beyond MMLU: Generalization to Broader Capabilities

Next, we ask how the proposed knowledge-infusion recipe in ACER transfers beyond MMLU. To this end, we evaluate the knowledge-infused Llama 3.2 models against their pretrained baselines on standard language understanding benchmarks. Table 2 reports results on ARC, GPQA, AGIEval, GSM8K, and HellaSwag. Tasks that emphasize domain-specific knowledge and recall, such as

ARC, GPQA, and MMLU (5-shot), benefit the most. Both the 3B and 1B variants achieve gains of more than two absolute percentage points, highlighting the effectiveness of targeted continual pre-training in strengthening reasoning over underrepresented knowledge areas.

Equally important, these improvements do not come at the cost of general capabilities. On benchmarks such as AGIEval, GSM8K, and HellaSwag, which measure general reasoning, arithmetic, and commonsense capabilities, the ACER models remain stable relative to their pretrained baselines. For the 3B model, differences are within 0.3 absolute points on HellaSwag and even smaller across other tasks, while the 1B model shows comparable stability. This indicates that domain infusion not only enhances specialized competence but also preserves broad language abilities, showing the scalability of our framework across diverse evaluation settings.

5 Conclusion

We present ACER (Automated Curriculum-Enhanced Regimen), a novel framework that systematically infuses domain-specific knowledge into large language models through progressive, curriculum-aligned learning sequences interspersed with targeted assessment. Our empirical evaluation demonstrates that ACER achieves substantial performance gains in target domains while maintaining model capabilities across existing tasks. This work establishes a principled foundation for structured knowledge integration in LLMs, offering a scalable pathway toward more capable and domain-aware language models.

646 **Limitations**

647 While ACER demonstrates consistent gains across
648 multiple domains and the tested model scales, sev-
649 eral limitations remain. First, although we evaluate
650 sensitivity to different synthetic data generators,
651 the quality and structural properties of the gener-
652 ated curriculum can influence the magnitude of
653 achievable gains (see Appendix D.2 for a qualita-
654 tive analysis). Designing synthesis pipelines that
655 explicitly optimize for curriculum structure and
656 granularity remains an open direction. Second, our
657 experiments focus on 3B and 1B parameter models.
658 While this allows controlled analysis of domain
659 knowledge infusion and curriculum effects, it re-
660 mains to be explored how ACER scales to larger
661 models and whether similar curriculum-aligned
662 training yields comparable or amplified benefits
663 at higher capacities.

664 **Ethical considerations**

665 This work focuses on improving domain knowl-
666 edge acquisition in large language models through
667 synthetic data generation and curriculum-aligned
668 continual pretraining. Our experiments use only
669 publicly available benchmarks and automatically
670 generated data, and do not involve human subjects
671 or personal data. As such, we do not foresee direct
672 ethical risks arising from this work. However, as
673 with all research on language models, downstream
674 use may inherit broader societal risks related to bi-
675 ased or unintended model behavior. We view ACER
676 as a training-time framework, and emphasize that
677 responsible deployment, evaluation, and monitor-
678 ing of models remain essential in practical settings.

679 **References**

680 Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien
681 Bubeck, Ronen Eldan, Suriya Gunasekar, Michael
682 Harrison, Russell J Hewett, Mojan Javaheripi, Piero
683 Kauffmann, and 1 others. 2024. Phi-4 technical re-
684 port. *arXiv preprint arXiv:2412.08905*.

685 Vinayak Arannil, Neha Narwal, Sourav Sanjukta
686 Bhabesh, Sai Nikhil Thirandas, Darren Yow-Bang
687 Wang, Graham Horwood, Alex Anto Chirayath, and
688 Gouri Pandeshwar. 2024. Dopamine: Domain-
689 specific pre-training adaptation from seed-guided
690 data mining. *arXiv preprint arXiv:2410.00260*.

691 Loubna Ben Allal, Anton Lozhkov, Guilherme Penedo,
692 Thomas Wolf, and Leandro von Werra. 2024. *Cos-*
693 *mopedia*.

Yoshua Bengio, Jérôme Louradour, Ronan Collobert, 694
and Jason Weston. 2009. *Curriculum learning*. In 695
Proceedings of the 26th Annual International Confer- 696
ence on Machine Learning, ICML '09, page 41–48, 697
New York, NY, USA. Association for Computing 698
Machinery. 699

Louis Béthune, David Grangier, Dan Busbridge, 700
Eleonora Gualdoni, marco cuturi, and Pierre Ablin. 701
2025. *Scaling laws for forgetting during finetuning* 702
with pretraining data injection. In *Forty-second In-* 703
ternational Conference on Machine Learning. 704

B. S. Bloom, M. B. Engelhart, E. J. Furst, W. H. Hill, 705
and D. R. Krathwohl. 1956. *Taxonomy of educational* 706
objectives. The classification of educational goals. 707
Handbook 1: Cognitive domain. Longmans Green, 708
New York. 709

Zhou Chen, Ming Lin, Zimeng Wang, Mingrun Zang, 710
and Yuqi Bai. 2024. Preparedllm: effective pre- 711
pretraining framework for domain-specific large lan- 712
guage models. *Big Earth Data*, 8(4):649–672. 713

Daixuan Cheng, Yuxian Gu, Shaohan Huang, Junyu Bi, 714
Minlie Huang, and Furu Wei. 2024. *Instruction pre-* 715
training: Language models are supervised multitask 716
learners. *Preprint*, arXiv:2406.14491. 717

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, 718
Ashish Sabharwal, Carissa Schoenick, and Oyvind 719
Tafjord. 2018. *Think you have solved question* 720
answering? try arc, the ai2 reasoning challenge. 721
Preprint, arXiv:1803.05457. 722

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, 723
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias 724
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro 725
Nakano, Christopher Hesse, and John Schulman. 726
2021. Training verifiers to solve math word prob- 727
lems. *arXiv preprint arXiv:2110.14168*. 728

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, 729
Abhinav Pandey, Abhishek Kadian, Ahmad Al- 730
Dahle, Aiesha Letman, Akhil Mathur, Alan Schel- 731
ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh 732
Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi- 733
tra, Archie Sravankumar, Artem Korenev, Arthur 734
Hinsvark, and 542 others. 2024. *The llama 3 herd of* 735
models. *Preprint*, arXiv:2407.21783. 736

Suchin Gururangan, Ana Marasović, Swabha 737
Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, 738
and Noah A. Smith. 2020. *Don't stop pretraining:* 739
Adapt language models to domains and tasks. In 740
Proceedings of the 58th Annual Meeting of the 741
Association for Computational Linguistics, pages 742
8342–8360, Online. Association for Computational 743
Linguistics. 744

Yao He, Xuanbing Zhu, Donghan Li, and Hongyu 745
Wang. 2025. *Enhancing large language models for* 746
specialized domains: A two-stage framework with 747
parameter-sensitive lora fine-tuning and chain-of- 748
thought rag. *Electronics*, 14(10). 749

750	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding . In <i>International Conference on Learning Representations</i> .	806
751		807
752		808
753		809
754		810
755	Jianheng Huang, Leyang Cui, Ante Wang, Chengyi Yang, Xinting Liao, Linfeng Song, Junfeng Yao, and Jinsong Su. 2024. Mitigating catastrophic forgetting in large language models with self-synthesized rehearsal . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1416–1428, Bangkok, Thailand. Association for Computational Linguistics.	811
756		812
757		813
758		814
759		815
760		
761		816
762		817
763	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models . <i>CoRR</i> , abs/2001.08361.	818
764		819
765		
766		820
767		821
768	Tobias Kerner. 2024. Domain-specific pretraining of language models: A comparative study in the medical field . <i>arXiv preprint arXiv:2407.14076</i> .	822
769		823
770		
771	Bruce W Lee, Hyunsoo Cho, and Kang Min Yoo. 2024. Instruction tuning with human curriculum . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 1281–1309, Mexico City, Mexico. Association for Computational Linguistics.	824
772		825
773		826
774		827
775		828
776		829
777	Haoran Li, Qingxiu Dong, Zhengyang Tang, Chaojun Wang, Xingxing Zhang, Haoyang Huang, Shaohan Huang, Xiaolong Huang, Zeqiang Huang, Dongdong Zhang, Yuxian Gu, Xin Cheng, Xun Wang, Si-Qing Chen, Li Dong, Wei Lu, Zhifang Sui, Benyou Wang, Wai Lam, and Furu Wei. 2025. Synthetic data (almost) from scratch: Generalized instruction tuning for language models . <i>Transactions on Machine Learning Research</i> .	830
778		831
779		832
780		833
781		834
782		
783		835
784		836
785	Huu Tan Mai, Cuong Xuan Chu, and Heiko Paulheim. 2024. Do llms really adapt to domains? an ontology learning perspective . In <i>International Workshop on the Semantic Web</i> .	837
786		838
787		839
788		840
789	Sven Najem-Meyer, Frédéric Kaplan, and Matteo Romanello. 2025. Don't stop pretraining! efficiently building specialised language models in resource-constrained settings . In <i>Proceedings of the 9th Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature (LaTeCH-CLfL 2025)</i> , pages 252–260, Albuquerque, New Mexico. Association for Computational Linguistics.	841
790		842
791		
792		843
793		844
794		845
795		846
796		847
797		848
798		849
799		850
800	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Driani, Julian Michael, and Samuel R. Bowman. 2023. Gpqa: A graduate-level google-proof qa benchmark . <i>Preprint</i> , arXiv:2311.12022.	851
801		
802		852
803		
804	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. Self-instruct: Aligning language models with self-generated instructions . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13484–13508, Toronto, Canada. Association for Computational Linguistics.	853
805		854
		855
		856
		857
		858
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858

Persona	Microeconomics	Statistics	Econometrics	Mathematics	Psychology
Highschool	0.67	0.74	0.55	1.04	0.72
Undergraduate	1.56	1.47	1.67	1.97	1.01
Graduate	2.56	2.10	2.73	1.55	1.27
Researcher	1.66	1.79	2.83	2.21	1.85
Total	6.46	6.10	7.78	6.77	4.85

Table 3: ACER book corpus: #tokens (in Millions) across the five domains and four personas. Total synthetic data #tokens = 31.97 Millions

Model	MEco _{hs}	Stats _{hs}	Econ	Maths _{hs}	Psych _p	Macro _t	Macro _{nt}
Cosmo-1B	0.2647	0.213	0.2719	0.2556	0.2663	0.2543	0.2626
Phi-1.5	0.4328	0.2685	0.2719	0.37	0.3971	0.3229	0.4266
Phi-2	0.5882	0.4259	0.3158	0.2926	0.5392	0.4323	0.5552
Llama 3.2 1B	0.3361	0.25	0.193	0.2296	0.3448	0.2707	0.4011
Llama 3.2 3B	0.5378	0.3796	0.3158	0.2704	0.5507	0.4108	0.5754

Table 4: Contextual comparison with synthetic-pretrained models: Cosmo 1B (1B parameters), Phi 1.5 (1.3B parameters), and Phi 2 (2.7B parameters), evaluated on the same MMLU subdomains used in this work. Results are provided for contextual reference and are not directly comparable to ACER-trained models due to differences in architecture and training setup.

- 859 2. Enables curriculum-aware scaling, where sub-
860 sequent synthetic data creation (e.g., text-
861 books and assessments) follows a structured
862 hierarchy rather than ad-hoc sampling.

863 A.2 Outline Generation

864 ACER leverages a large language model prompted
865 as a professional author to create a **multilevel out-**
866 **line** in a structured JSON format. The generated
867 outline includes the following.

- 868 • **Book Title:** A descriptive title aligned with
869 the topic and purpose.
- 870 • **Parts:** Four to six major parts that divide the
871 book into thematic areas, each with a concise
872 summary.
- 873 • **Chapters:** Four to six chapters per part,
874 each providing a self-contained exploration
875 of subtopics.
- 876 • **Sections and Subsections:** Three to six sec-
877 tions per chapter, with optional subsections
878 for complex concepts, ensuring fine-grained
879 coverage.

880 This structured representation balances clarity,
881 depth, and scalability, allowing the same frame-
882 work to create outlines of varying complexity based
883 on audience needs.

Explicit outline generation provides three main
benefits:

- 886 1. **Consistency:** Maintains structural uniformity
887 across different domains and audience levels.
- 888 2. **Scalability:** Simplifies automation, as the
889 JSON schema can be directly consumed by
890 subsequent content generation pipelines.
- 891 3. **Pedagogical Alignment:** Encourages system-
892 atic progression from foundational concepts
893 to advanced material, supporting curriculum-
894 driven pretraining.

895 A.3 Synthetic Content Creation

896 The procedure begins with the construction of a
897 tree representation from the ToC. Each node in this
898 tree corresponds to a textual unit—either a *Part*,
899 *Chapter*, *Section*, or *Subsection*. For every node,
900 the system maintains key attributes that include
901 the title, description, content, and node type (an
902 enumerated label designating whether the node be-
903 longs to root, part, chapter, section, or subsection).
904 This representation provides a flexible yet struc-
905 tured foundation for downstream content synthesis.

906 Once the hierarchical structure is established,
907 the system proceeds to parse the tree and generate
908 section-level content. For each node in the section,
909 the input prompt is carefully composed to include

Model (Llama 3.2 1B)	MEco _{hs}	Stats _{hs}	Econ	Maths _{hs}	Psych _p	Macro _t	Macro _{nt}
Pretrained baseline	0.3361	0.25	0.193	0.2296	0.3448	0.2707	0.4011
Cog+Con _{gemini-2.0-flash}	0.3445	0.3287	0.2105	0.2481	0.3562	0.2976	0.4051
Cog+Con _{Phi-4}	0.3361	0.2917	0.2105	0.2593	0.3317	0.2859	0.4017
Cog+Con _{Llama 3.1 70B}	0.3529	0.3241	0.2193	0.2407	0.3284	0.2931	0.3962

Table 5: Comparison of pretrained Llama 3.2 1B with ACER-trained variants under the best-performing curriculum (Cog+Con), using different synthetic data generators. Domain-targeted synthetic textbooks with exam-style QA improve performance over the pretrained baseline, with Cog+Con (Gemini 2.0 Flash) achieving a **2.7**-point improvement on target domains. Results across Phi-4 and Llama 3.1 70B generators show that while generator choice modulates the magnitude of gains, curriculum-aligned training improves or preserves performance without degradation.

contextual information such as the title and description of the enclosing part, the title and description of the chapter, the focal section title, and the list of subsections under it. In addition, metadata style guidelines (intent, audience, genre, tone, voice, and language) are explicitly injected into the prompt. This ensures that the generation process is not only content-driven but also stylistically aligned with the pedagogical and editorial goals of the textbook. The LLM is instructed to produce polished, instructional prose while avoiding any meta-commentary, annotations, or explanations of its reasoning.

B Tokens generated using ACER synthesis flow

Table 3 details the number of tokens generated by ACER. It uses Gemini 2.0 Flash apis for synthetic data generation.

C Positioning Relative to Synthetic Pretraining Models

We contextualize our work relative to recent pre-training efforts that expand knowledge coverage via large-scale synthetic data. We evaluate representative 1B-3B parameter models trained predominantly on broad synthetic corpora, specifically Phi 1.5, Phi 2, and Cosmopedia, on the same five MMLU subdomain considered in this work. These models reflect alternative strategies that rely on broad, general-purpose synthetic data rather than targeted curriculum-aligned knowledge infusion.

As reported in Table 4, these models exhibit low performance in specialized domains such as economics and professional psychology, indicating limited domain-specific knowledge despite their scale and diversity. We emphasize that these results are not directly comparable to ACER-trained models due to differences in model architecture,

scale, and training objectives. Nevertheless, this analysis highlights that broad synthetic pretraining alone does not reliably close fine-grained domain gaps, motivating the need for targeted, curriculum-aware knowledge infusion as explored in ACER.

D Experiments: Additional Details

D.1 Reproducibility Details

The parameters required to reproduce our results, such as token budgets, sequence lengths, batch sizes, optimizer settings, learning rate schedule, training steps, hardware specifications, random seeds, prompts, and de-duplication thresholds—are provided in section 4.1. The prompts used for all types of generations are listed in Appendix G.

D.2 Robustness to Synthetic Data Generators

To assess whether ACER’s gains depend on a specific synthesis model, we regenerate the synthetic textbook corpus using Phi-4 and Llama 3.1 70B, while fixing the curriculum to the best-performing variant (Cog+Con) and keeping all other settings unchanged. As shown in Table 1 (bottom block), Cog+Con with Phi-4 generated data continues to improve over the pretrained Llama 3.2 3B baseline on target MMLU domains, whereas Llama 3.1 70B generated data yields smaller gains and performs comparably to the baseline. These results show that ACER is robust to the choice of synthetic data generator: even when using weaker or less structurally detailed generators (e.g. Llama 3.1 70B) instead of Gemini 2.0 Flash, training does not lead to performance degradation, although stronger generators yield larger gains. We find that this behavior correlates with the structural granularity of the synthesized textbooks: Gemini-generated books contain finer-grained outlines (about 100 sections per book on average) compared to Llama 3.1 70B

Ratio	Baseline	1:3	1:1	3:1	9:1
Target Macro	0.4108	0.4257	0.4407	0.4284	0.4135
Non-Target Macro	0.5754	0.5759	0.5821	0.5805	0.5736

Table 6: Ablation on in-domain to replay data proportion for training: CPT with equal proportion of in-domain and replay data provides the right balance

(about 30 sections). Since ACER’s curricula operate at the section level, this reduced granularity weakens the progressive supervision signal during continual pretraining. Corresponding results for Llama 3.2 1B follow similar trends and are reported in Appendix D.3.

D.3 Results: Scaling to Llama 3.2 1B

We apply the best-performing Cog+Con curriculum to the Llama 3.2 1B model to assess whether ACER’s benefits extend to smaller model scales. As shown in Table 5, the Cog+Con variant achieves improvements of up to 2.7 points on target MMLU domains relative to the pretrained baseline, while largely preserving performance on non-target domains.

We further evaluate robustness to the choice of synthetic data generator by regenerating the synthetic corpus using Phi-4 and Llama 3.1 70B, while keeping the curriculum schedule fixed. Consistent with observations at the 3B scale, Phi-4 generated data yields comparable improvements, whereas Llama 3.1 70B generated data leads to smaller or neutral gains. These results suggest that, even at smaller scales, ACER does not rely on a specific generator to avoid degradation, although generator choice influences the magnitude of achievable gains.

D.4 Ablations: What is the Right Data Mixture?

Table 6 reports the effect of varying the ratio of in-domain synthetic data to replay data on model performance under the Cog+Con curriculum. The balanced 1:1 setting yields the strongest results, improving Macro_t by about 3 absolute points over the baseline and Macro_{nt} by 0.7 points. Ratios skewed toward replay (1:3) or synthetic data (3:1, 9:1) yield smaller gains and fail to surpass the balanced mixture. The trend is particularly evident for non-target domains: as replay proportion decreases, Macro_{nt} consistently drops, underscoring the role of replay in preserving generalization. Conversely, in-domain synthetic data is the main driver

of improvements on target domains. Overall, the balanced 1:1 mixture offers the best trade-off between specialization and generalization, pointing to adaptive mixture schedules as a promising future direction.

E Example domain: Microeconomics

This section shows the generated book content for **Microeconomics** domain (persona: high school), providing snippets for the following:

- Topic detailing 1033
- ToC (shortened version) 1034
- Example section 1035
- Example QA Pair 1036
- Example Easy QA Pair 1037
- Example Hard QA Pair 1038
- A snapshot of the ToC Tree Visualiser 1039

Topic Detailing: Imagine you’re running a lemonade stand. Microeconomics is like understanding all the tiny decisions that make your stand successful 2013 how many lemons to buy, what price to charge for your lemonade, and whether to hire a friend to help you out. It’s about how individuals, families, and businesses make choices about using limited resources. Instead of looking at the entire economy of a country (that’s macroeconomics!), microeconomics zooms in on these smaller, everyday economic activities. It helps us understand why things cost what they do, how markets work, and how people respond to different incentives. Think of it as the science of small-scale economic decisions and their ripple effects.

1: Shortened ToC (JSON)

```
{
  "title": "Part 1: The Basics of Making Choices",
  "description": "This section introduces the foundational.... ",
  "chapters": [

```

```

1062 {
1063   "title": "Chapter 1: Welcome to
1064     Microeconomics!",
1065   "description": "This chapter
1066     introduces the world of
1067     microeconomics ...",
1068   "sections": [
1069     {
1070       "title": "What is
1071         Microeconomics?",
1072       "subsections": [
1073         {
1074           "title": "The 'Micro' in
1075             Microeconomics"
1076         },
1077         {
1078           "title": "Lemonade Stands
1079             and Economies"
1080         }
1081       ]
1082     },
1083     {
1084       "title": "Why Should You Care
1085         About Microeconomics?",
1086       "subsections": [
1087         {
1088           "title": "Making Smarter
1089             Choices"
1090         },
1091         {
1092           "title": "Understanding
1093             the World Around You"
1094         }
1095       ]
1096     },
1097     :
1098     .
1099   ]
1100 }
1101 ]
1103 }

```

1104 **Example Section:**
1105 **Part title:** The Basics of Making Choices
1106 **Chapter title:** Scarcity: The Reason We Have to
1107 Choose
1108 **Section title:** Dealing With Scarcity: Prioritization
1109 and Rationing
1110 **Section content:** Scarcity, as we've established,
1111 means we can't have everything we want. This
1112 universal truth forces us to confront some tough
1113 questions. If we can't satisfy all our desires, how
1114 do we decide which ones to satisfy and which ones
1115 to leave unfulfilled? The answer lies in two closely
1116 related concepts: prioritization and rationing. ...

1117 **Example QA :**
1118 **Question:** What are the three main decision-
1119 makers studied in microeconomics?
1120 **Answer:** Microeconomics focuses on the "small
1121 picture" of the economy, studying the decisions
1122 made by individual players rather than the entire

national or global economy..... 1123

Example Easy QA : 1124

Question: Imagine you're deciding whether to
1125 spend your allowance on a new video game or save
1126 it for a concert ticket. How does this single deci-
1127 sion illustrate the core idea that microeconomics
1128 studies choices made by individuals with limited
1129 resources? 1130

Answer: Okay, that's a great question! It gets right
1131 to the heart of what microeconomics is all about.
1132 Let's break it down using the lemonade stand ex-
1133 ample and the ideas presented in the text..... 1134

Example Hard QA : 1135

Question: Imagine two lemonade stands on the
1136 same street. One stand uses only organic lemons
1137 and charges a higher price, while the other uses
1138 regular lemons and charges a lower price..... 1139

Answer: Okay, that's a great question! It gets
1140 right to the heart of how microeconomics works
1141 in the real world, even in something as simple as
1142 two lemonade stands. Microeconomics helps us
1143 understand why both lemonade stands 2013 the one
1144 with organic lemons and a higher price and the.... 1145

Tree Visualiser Tool : Figure 4 shows the tree vi-
1146 sualizer tool we developed to analyze the generated
1147 data. 1148

**F Example of Performance Improvement
after ACER** 1149 1150

We present an example MMLU question from the
1151 high-school microeconomics subset to demonstrate
1152 the effectiveness of ACER. The base pre-trained
1153 model (LLaMA 3.2 3B) could not answer this ques-
1154 tion correctly, while the ACER trained model was
1155 able to find the right answer due to its enhanced
1156 domain understanding. 1157

Question: "The elasticity of supply is
1158 typically greater when ..."
1159

Options: 1160

- A. Producers have fewer alternative goods to pro-
1161 duce. 1162
- B. Producers have less time to respond to price
1163 changes. 1164
- C. Producers are operating near the limits of their
1165 production. 1166
- D. Producers have more time to respond to price
1167 changes. 1168

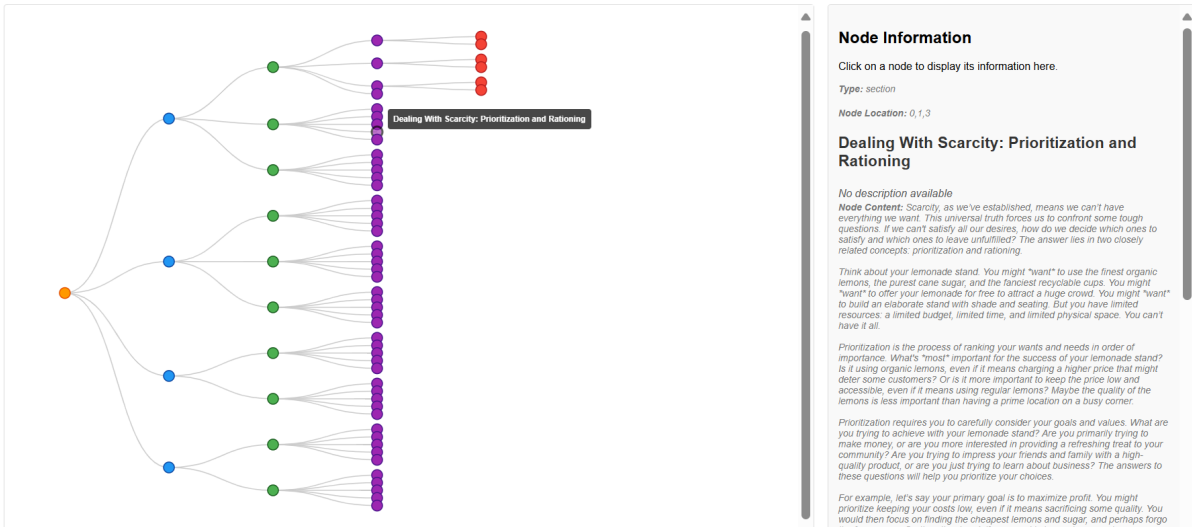


Figure 4: Snapshot of the Tree Visualisation tool we created to analyse the ToC and the associated generated Data

The correct answer is option D. The baseline LLaMA 3.2 3B model predicted option A (incorrect), while our ACER trained model predicted option D (correct).

This improvement can be attributed to the additional knowledge injected during continual pretraining. Our generated microeconomics book included a dedicated section on “Supply Curve and Elasticity” An excerpt from the Table of Contents of that book is shown below:

2: Excerpt from the Table of Contents (JSON)

```
{
  "title": "Advanced Producer Theory: Cost and Production",
  "description": "This chapter focuses on the theory of the firm...",
  "sections": [
    {
      "title": "Profit Maximization and Supply",
      "subsections": [
        {
          "title": "Supply Curve and Elasticity"
        }
      ]
    },
    {
      "title": "Applications of Producer Theory",
      "subsections": [
        {
          "title": "Impact of Technological Change"
        }
      ]
    },
    {
      "title": "Summary"
    }
  ]
}
```

}

1219

G Prompts used in the ACER framework

1212

G.1 Domain Detailing Prompt

1213

The prompt used for domain detailing:

1214

3: Prompt for Domain Detailing

```
prompt = f"""
You are an expert author preparing to write a comprehensive book.
Topic: \"{topic}\".
Target audience: \"{audience}\".
Intent (purpose): \"{intent}\".

Your task is to generate preparatory material that will help you structure the book. Based on the above, provide the following:
1. A factual and high-level description of the topic, suitable for the target audience.
2. A list of 6-8 core themes or subtopics that must be covered to provide a well-rounded understanding of the topic.
3. A list of 6-8 important questions that the book should aim to answer from the perspective of the target audience.

Present your output in JSON with the following keys:
- "description": A string
- "subtopics": A list of strings
- "key_questions": A list of strings
"""
```

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1248

G.2 Outline Generation Prompt

1249

The prompt used for outline (Table of Contents) generation:

1250

4: Prompt for Outline generation

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

```

prompt = f'''
You are a professional book author known
for clear, structured, and reader-
friendly writing.
Your task is to design a full book
outline based on the information
below.
---
**Topic:** {topic}

**Target Audience:** {audience}

**Intent:** {intent}

**Genre:** {genre}

**Tone:** {tone}
**Voice:** {voice}
**Language Style:** {language}

**Topic Description:**
{description}

**Core Subtopics to Cover:**
{subtopics}

**Key Questions the Book Should Answer
:**
{key_questions}
---

Based on this, generate a comprehensive
book structure in JSON with the
following:
1. A compelling and descriptive "title"
for the book.
2. Divide the book into **4 to 6 major
parts**:
- Each part should have:
  - "title" (clear, theme-based)
  - "description" (3 5 sentence
summary)
3. Each part should contain **4 to 6
chapters**:
- Each chapter should include:
  - "title"
  - "description" (3 5 sentences,
aligned with purpose)
  - "sections": 3 to 6 entries, each
with:
    - "title" (section title)
    - "subsections": Optional (to
describe complex ideas, 2-3
subsections per section)
    - "title" (subsection title
)
(Include common sections like "
Introduction" and "Summary" where
appropriate.)

Return **only** a well-formatted JSON
object in the following structure:

```

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

1465

1466

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

1512

1513

1514

1515

1516

1517

1518

1519

1520

1521

1522

1523

1524

1525

1526

1527

1528

1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

```

1382 - Ensure thematic and stylistic
1383     continuity with the previous content
1384     .
1385 - Do not include headers, instructions,
1386     or artificial markers only write
1387     the finished prose.
1388 ...

```

1390 **G.4 QA Pair Generation Prompts**
1391 The prompt used for question generation is:

1392 6: Prompt for question generation

```

1393 prompt = f"""
1394 You are a professor specializing in the
1395 subject of {topic_name}. Your task
1396 is to generate **educational, self-
1397 contained questions** to help
1398 students understand a section from a
1399 textbook titled: "{book_title}".
1400 This book is for {target_audience}, with
1401 the primary goal to: **{intent}**
1402
1403 ### Global Context Start ###
1404 - **Topic**: {topic_name}
1405 - **Topic Description**: {
1406     topic_description}
1407 - **Questions the book attempts to
1408     answer**: {guiding_questions}
1409 ### Global Context End ###
1410
1411 You will be generating questions for the
1412 following part of the book:
1413 - Chapter: {chapter_title}
1414 - Section: {section_title}
1415 --- Section Content Start ---
1416 {section_text}
1417 --- Section Content End ---
1418 """
1419
1420 if previous_question:
1421     prompt += f"""
1422 ### TASK ###
1423
1424 You are generating the next question in
1425 a sequence.
1426
1427 1. Analyze the difficulty of the
1428     previous question.
1429 2. Generate a new question that is **
1430     slightly more challenging** than the
1431     one before. Increase cognitive
1432     depth move progressively from
1433     understanding interpretation
1434     analysis application.
1435 3. Ensure the question is different in
1436     focus or angle from the previous one
1437     .
1438 4. Do **not** include the answer.
1439 5. **Strictly** maintain relevance to
1440     the given section and suitability
1441     for {target_audience}.
1442 6. **Only output the question** do not
1443     include any prefixes like "
1444     Generated Question", "Q1:", or
1445     anything else.
1446
1447 Previous Question:
1448 - {previous_question}

```

```

1449 ### Question ###
1450 Question: """
1451
1452 else:
1453     prompt += f"""
1454 ### TASK ###
1455
1456 You are generating the **first question
1457 ** in a learning sequence.
1458
1459 1. Start with a **simple question**
1460     focused on factual recall or basic
1461     definitions.
1462 2. Make the question self-contained and
1463     directly based on the section
1464     content.
1465 3. Do **not** include the answer.
1466 4. Ensure clarity, relevance, and
1467     suitability for {target_audience}.
1468 5. **Only output the question** do not
1469     include any prefixes like "
1470     Generated Question", "Q1:", or
1471     anything else.
1472
1473 ### Question ###
1474 Question: """

```

The prompt used for answer generation is:

7: Prompt for answer generation

```

1477 prompt = f"""
1478 You are a professor specializing in the
1479 subject of {topic_name}. Your task
1480 is to generate an **educational,
1481 self-contained answer** to a
1482 question based on a section from a
1483 textbook titled: "{book_title}".
1484 This book is for {target_audience}, with
1485 the primary goal to: **{intent}**
1486
1487 ### Global Context Start ###
1488 - **Topic**: {topic_name}
1489 - **Topic Description**: {
1490     topic_description}
1491 ### Global Context End ###
1492
1493 You are answering a question based on
1494 the following section:
1495 - Chapter: {chapter_title}
1496 - Section: {section_title}
1497
1498 --- Section Content Start ---
1499 {section_text}
1500 --- Section Content End ---
1501
1502 ### TASK ###
1503
1504 Your job is to write a complete,
1505 thoughtful, and educational answer
1506 to the student's question below.
1507 Follow these guidelines:
1508
1509 1. Ensure the answer is **directly
1510     grounded in the section content**
1511     provided above.
1512 2. Make the answer **rich in detail**,
1513     using clear language and examples
1514     when helpful.
1515 3. The answer should be **self-contained
1516

```

```

1517     **, meaning the reader should not
1518     need to refer to the original
1519     content to understand it.
1520 4. Ensure the tone is educational and
1521     appropriate for {target_audience}.
1522 5. Do **not invent any information** not
1523     present in the section.
1524
1525 ### Question ###
1526 {question}
1527
1528 ### Answer ###
1529 """
1530

```

1531 H Decontamination- data analysis

1532 We performed semantic deduplication between the
1533 synthetic data we generated and the benchmark
1534 datasets. The motivation for using semantic dedu-
1535 plication is that it is substantially more thorough
1536 than MinHash-based deduplication, as it can detect
1537 near-duplicates beyond surface-level lexical over-
1538 lap. Specifically, we first computed embeddings
1539 for both the generated samples and the benchmark
1540 data. For each generated sample, we then identified
1541 its nearest neighbors in the benchmark set using
1542 cosine similarity. Samples with similarity scores
1543 above a predefined threshold were manually ana-
1544 lyzed. Across all domains, we observed little to
1545 no contamination with respect to the benchmark
1546 data. For the few cases where the cosine similarity
1547 exceeded 0.9, the generated samples were found
1548 to discuss canonical concepts that are commonly
1549 covered when acquiring foundational knowledge
1550 of the domain, rather than reproducing benchmark-
1551 specific content. A detailed decontamination report
1552 is provided in Table 7, along with representative
1553 example fragments in Table 8.

Domain	Generated Samples	Similarity ≥ 0.8
Econometrics	4125	0.0969
Mathematics	4110	0.0243
Psychology	4350	0.2758
Statistics	4200	0.1428
Microeconomics	3900	0.051282

Table 7: Domain statistics for generated data (similarity threshold = 0.8). The final column reports the percentage of benchmark samples with cosine similarity ≥ 0.8

Contents
<p>Domain: Econometrics</p> <p>Generated Sample: What is the primary consequence of omitting a relevant variable from a multiple linear regression model?</p> <p>Closest Benchmark Sample: If a relevant variable is omitted from a regression equation, the consequences would be that:</p> <p>i) The standard errors would be biased ii) If the excluded variable is uncorrelated with all of the included variables, all of the slope coefficients will be inconsistent. iii) If the excluded variable is uncorrelated with all of the included variables, the intercept coefficient will be inconsistent. iv) If the excluded variable is uncorrelated with all of the included variables, all of the slope and intercept coefficients will be consistent and unbiased but inefficient.</p> <p>Cosine Similarity: 0.8456</p>
<p>Domain: Mathematics</p> <p>Generated Sample: A university club with 15 members needs to form a committee of 5 to organize a fundraising event. * How many different committees can be formed? * The club decides that the committee must have a president, a vice-president, a treasurer, a secretary, and a public relations officer. The president and vice-president must be chosen from among 8 senior members, while the treasurer, secretary, and public relations officer can be chosen from the remaining members (seniors or otherwise). In how many ways can such a committee be formed, ensuring each member has a distinct role?</p> <p>Closest Benchmark Sample: How many different possible committees of 5 people can be chosen from a group of 15 people?</p> <p>Cosine Similarity: 0.8338</p>
<p>Domain: Psychology</p> <p>Generated Sample: What is the primary difference between classical conditioning and operant conditioning?</p> <p>Closest Benchmark Sample: What is the major difference between classical and operant conditioning?</p> <p>Cosine Similarity: 0.9871</p>
<p>Domain: Statistics</p> <p>Generated Sample: What is a sampling distribution, and how is it constructed?</p> <p>Closest Benchmark Sample: What is a sampling distribution?</p> <p>Cosine Similarity: 0.929</p>
<p>Domain: Microeconomics</p> <p>Generated Sample: What fundamental economic problem does microeconomics primarily address?</p> <p>Closest Benchmark Sample: The primary focus of microeconomics is</p> <p>Cosine Similarity: 0.8312</p>

Table 8: Examples of generated samples, closest benchmark samples and their cosine similarity