

Tracing Multilingual Knowledge Acquisition Dynamics in Domain Adaptation: A Case Study of English-Japanese Biomedical Adaptation

Anonymous ACL submission

Abstract

Multilingual domain adaptation (ML-DA) is widely used to learn new domain knowledge across languages into large language models (LLMs). Although many methods have been proposed to improve domain adaptation, the mechanisms of multilingual knowledge acquisition, how domain knowledge is learned within a language and transferred across languages, remain underexplored. This gap leads to suboptimal performance, particularly in low-resource settings. This work examines the learning dynamics of LLMs during ML-DA. Because prior ML-DA studies often train and evaluate on datasets with mismatched knowledge coverage, we propose **AdaXEval**, an adaptive evaluation method that builds multiple-choice QA datasets from the same bilingual domain corpus used for training, thereby directly studying multilingual knowledge acquisition. Through continual training of LLMs with diverse data recipes, we track how LLMs acquire domain facts and pinpoint the mechanism behind the transformation process from domain training data to knowledge. Our experiments on a 13B English-Japanese bilingual LLM reveal that cross-lingual transfer remains challenging despite a high-quality bilingual corpus. The code has been released. [MLDA-Eval](#)

1 Introduction

Large language models (LLMs) trained on general-domain corpora perform well on diverse tasks but struggle in specialized domains (Jang et al., 2022b,a). Domain adaptation addresses this by continually training LLMs on domain-specific data to enhance expertise (Jiang et al., 2024a; Lai-king and Paroubek, 2024; Xie et al., 2024b). Although prior work has explored strategies such as data augmentation and cross-lingual transfer to improve adaptation efficiency on low-resource settings (Fang et al., 2023; Gao et al., 2024a), the mechanisms underlying effective domain knowledge acquisition and transfer remain insufficiently understood.

Understanding domain adaptation requires examining how LLMs acquire and internalize facts from domain data. Prior studies show that factual knowledge accumulates gradually through repeated exposures, shaped by data frequency and model scale (Chang et al., 2024; Liu et al., 2025; Zhao et al., 2024b). Moreover, multilingual settings introduce additional complexity, as equivalent facts may be encoded differently across languages (Mondal et al., 2025; Zhao et al., 2024b). However, Most analyses focus on predefined relational probes rather than real domain facts with complex structures and specialized terminology. Moreover, the link between training data and acquired knowledge also remains underexplored, which is crucial to optimizing domain adaptation strategies.

Our work aims to investigate the process of knowledge acquisition in domain adaptation from a mechanistic perspective. Specifically, we seek to understand, during the continual-training process on the domain corpus, how domain facts are **memorized** and **generalized** across different linguistic contexts, including both **intralingual** (within a language) and **interlingual** (across languages) variations, and to identify the key factors that facilitate effective knowledge acquisition and transfer. To achieve the goal, we focus on three questions:

- RQ1:** How to effectively evaluate the domain knowledge acquisition from diverse aspects?
- RQ2:** What is the mechanism behind the transformation from training data to knowledge?
- RQ3:** What factors are critical to achieve cross-lingual transfer?

Existing approaches to evaluate domain knowledge primarily rely on public benchmarks (Singhal et al., 2022; Jiang et al., 2025) or training loss analysis (Zucchet et al., 2025; Liu et al., 2025). However, such benchmarks offer limited coverage for low-resource settings and fail to capture knowledge generalization abilities. Moreover, the

084 misalignment between training data and bench-
085 mark knowledge coverage makes evaluation an
086 imperfect reflection of acquired knowledge. To
087 address these gaps and resolve **RQ1**, we propose
088 **AdaXEval**, an adaptive domain knowledge evalua-
089 tion data generation pipeline. AdaXEval automati-
090 cally constructs multiple-choice datasets to evalu-
091 ate knowledge *memorization*, intralingual general-
092 ization (*paraphrase*), and interlingual generaliza-
093 tion (*cross-lingual transfer*). AdaXEval operates
094 on either a monolingual or a bilingual domain cor-
095 pus, where the latter is required for cross-lingual
096 transfer evaluation, ensuring broad applicability
097 across rare domains and low-resource languages.
098 Human annotation from multiple perspectives con-
099 firms that AdaXEval provides reliable evaluation.

100 We next investigate how training data is dy-
101 namically transformed into knowledge (**RQ2**).
102 Specifically, we conduct a case study of Japanese
103 biomedicine domain adaptation using a 13B En-
104 glish/Japanese bilingual LLM (LLM-jp et al.,
105 2024), with English serving as a comparison and
106 source for knowledge transfer. We begin with
107 monolingual continual training on both English
108 and Japanese using the J-STAGE corpus, which
109 contains biomedical documents for both languages.
110 By evaluating training checkpoints with AdaXEval-
111 generated datasets, we observe a gradual knowl-
112 edge acquisition process for cloze queries and para-
113 phrases; however, LLM struggles to achieve cross-
114 lingual transfer. Further analysis reveals that knowl-
115 edge is acquired as losses of correct options are
116 shielded from rapid growth due to the model’s
117 overfitting to training data, which we term **loss**
118 **shielding**. This acquisition eventually plateaus as
119 training causes the model to overfit significantly to
120 the training data, resulting in a substantial increase
121 in loss across all options in evaluation instances.
122 Examining losses on perturbed training data reveals
123 that LLMs readily overfit to fixed token sequences
124 in the training data, even under minor noise.

125 Finally, we investigate key factors influenc-
126 ing cross-lingual transfer of domain knowledge
127 through **multilingual continual training** with di-
128 verse data recipes (**RQ3**). We focus on transla-
129 tion and romanization strategies to enhance trans-
130 fer. Our results show that cross-lingual token
131 overlap in related domains is crucial for effec-
132 tive knowledge transfer. Nonetheless, even with
133 high-quality alignment signals such as translations,
134 cross-lingual transfer remains challenging, under-
135 scoring the need for more effective methods.

2 Knowledge Acquisition Evaluation 136

To effectively evaluate domain adaptation in low-
137 resource scenarios, we propose AdaXEval, an adap-
138 tive pipeline for generating evaluation datasets. 139

2.1 AdaXEval Pipeline 140

AdaXEval is an adaptive evaluation pipeline that
141 evaluates domain knowledge acquisition by gener-
142 ating evaluation datasets directly from the training
143 corpus, ensuring evaluated facts stay aligned with
144 the training data. The pipeline includes four steps:
145 fact detection, question crafting, distractor genera-
146 tion and quality filtering, as illustrated in Figure 11. 147

2.1.1 Fact Detection 148

AdaXEval first detects sentences that may contain
149 domain facts through a two-step strategy: named-
150 entity-recognition (NER)-based sentence filter-
151 ing and multi-agent fact detection. First, domain-
152 specific NER tools and linguistic heuristics are em-
153 ployed to identify sentences in the training corpora
154 that contain multiple named entities. Next, we
155 design Chain-of-Thought (CoT) instructions to de-
156 tect sentences containing domain facts from the
157 filtered sentences, and extract triples in the format
158 $\langle \text{subject}, \text{relation}, \text{object} \rangle$ as the reference
159 for question crafting. Specifically, AdaXEval em-
160 ploys a multi-LLM agent for fact detection and
161 triple extraction, estimating the overall confidence
162 of the outputs, and adapting the top-confident ex-
163 traction result to improve evaluation reliability. 164

2.1.2 Queries Crafting 165

Given the factual sentence and referenced triple,
166 we first prompt the LLM to generate reliable do-
167 main factual triples (*e.g.*, $\langle \text{blood sugar level, can}$
168 $\text{be controlled by, insulin} \rangle$). As complex domain
169 knowledge cannot be easily formalized into named
170 entities or predefined relations, subjects and objects
171 are preferably named entities, though descriptive
172 phrases are acceptable. AdaXEval uses advanced
173 LLMs to generate diverse question-answer pairs
174 measuring three knowledge acquisition abilities. 175

1) Knowledge memorization uses a cloze prompt
176 with [BLANK] as the placeholder for the object (*e.g.*,
177 *Blood sugar level can be controlled by [BLANK].*).
178 Given the original sentence and refined triple as
179 reference, we prompt the LLM to generate a cloze
180 question that closely matches the original sentence
181 to assess memorization exclusively. 182

2) Intralingual generalization assesses LLMs’ 183

184 ability to acquire knowledge using linguistic ex- 232
185 pressions that vary from those in the training cor- 233
186 pus. We design CoT instructions to let LLMs para- 234
187 phrase the cloze queries into question-like style 235
188 questions where different vocabulary is encouraged 236
189 (e.g., Which substance helps manage glycemic 237
190 levels in the body?) 238

191 **3) Interlingual generalization** measures how 239
192 learned facts can be transferred across languages. 240
193 While translation is a strong candidate, translating 241
194 sentences that express domain knowledge is chal- 242
195 lenging due to the presence of specialized named 243
196 entities, terminology, and concepts that lack direct 244
197 equivalents across languages (Liang et al., 2024). 245
198 To address this, we adapt AdaXEval to a bilingual 246
199 domain corpus containing languages X and Y , us- 247
200 ing the paraphrased dataset from language X to 248
201 evaluate cross-lingual transfer capabilities in Y . 249

202 2.1.3 Distractor Generation 250

203 AdaXEval then generates three plausible yet incor- 251
204 rect answer options that remain topically related 252
205 but unambiguously wrong, while explicitly instruct- 253
206 ing the advanced LLM to avoid surface-level cues 254
207 such as sequence length. 255

208 2.1.4 Quality Filtering 256

209 Finally, AdaXEval uses the LLM to filter low- 257
210 quality multiple-choice QA instances that fail to 258
211 meet the requirements in § 2.1.1,2.1.2,2.1.3. 259

212 2.2 Evaluation metric 260

213 For each evaluation dataset, we follow Gao et al. 261
214 (2023) to compute the average cross-entropy loss 262
215 over the target tokens of possible answers and se- 263
216 lect the one with the highest generation possibility 264
217 as the final answer. Specifically, for loss calculation 265
218 of cloze queries, we use tokens before the [BLANK] 266
219 as context and compute loss on the following to- 267
220 kens. For paraphrases, we treat the question as 268
221 context and measure only the loss of answer to- 269
222 kens. We use prediction accuracy as the metric for 270
223 knowledge acquisition. See Appendix B.2 for the 271
224 mathematical formulation of the evaluation metric. 272

225 2.3 Experimental Setup 273

226 **Domain corpus:** Our study investigates biomed- 274
227 ical domain adaptation in English–Japanese as a 275
228 case study. Specifically, we utilize the J-STAGE, an 276
229 English-Japanese bilingual biomedical corpus (see 277
230 § 3.1), as the data source for both model training 278
231 and AdaXEval generation. 279

Details of Generation: We randomly sampled 232
10,000 bilingual documents to generate the evalua- 233
tion dataset. We split abstracts into sentences and 234
filter out sentences with fewer than two biomedical 235
entities. For fact detection of filtered sentences, 236
we use three open-source LLMs ¹ from different 237
families to assess the confidence of whether the sen- 238
tence is factual for each language. We then sum the 239
confidence scores across the three LLMs and retain 240
sentences with combined confidence scores greater 241
than 2 (maximum 3). Finally, we use GPT-4.1 to 242
generate cloze queries, paraphrases, and three dis- 243
tractors for each instance. See Appendices B.1 and 244
C for details of the generation process and statisti- 245
cal report of generated datasets. 246

Human Evaluation: To assess the quality of our 247
generated datasets, we conduct a comprehensive 248
human evaluation across four key components of 249
the knowledge extraction and question generation 250
pipeline, including triple extraction quality evalua- 251
tion, cloze prompt evaluation, paraphrased question 252
evaluation, and distractor quality evaluation. Over- 253
all, the evaluation result indicates that AdaXEval is 254
able to generate high-quality evaluation data, meet- 255
ing the requirements for assessing diverse knowl- 256
edge acquisition abilities. See Appendix C for 257
evaluation results and dataset examples, and Ap- 258
pendix G for the human evaluation guideline. 259

260 3 Tracing Knowledge Acquisition 260

261 This section examines the training dynamics of 261
domain adaptation and explores the mechanism 262
underlying the transformation from training data 263
to knowledge in the monolingual setting. We per- 264
form monolingual continual training on English 265
and Japanese using a bilingual biomedical dataset. 266

267 3.1 Experimental Setup 267

Data preparation: We utilize a subset of the 268
J-STAGE corpus, which comprises Japanese re- 269
search papers with some abstracts translated into 270
English.² Specifically, we select 614,444 Japanese 271
and 404,643 English biomedical documents, paired 272
one-to-one. These bilingual pairs provide source 273
data for AdaXEval generation. To strengthen do- 274
main adaptation and enable fine-grained analysis, 275

¹We employ open-source LLMs for local inference, as the large number of candidate sentences would otherwise incur substantial computational costs.

²Access to the dataset is restricted by the J-STAGE license, so it cannot be publicly released.

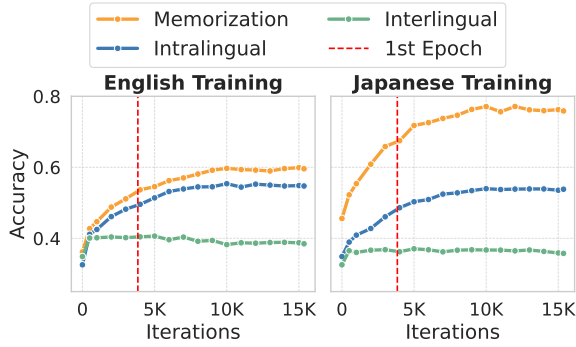


Figure 1: Dynamic knowledge acquisition evaluation after monolingual pretraining.

we apply instruction pretraining as a data augmentation baseline (Jiang et al., 2024b). Biomedical instructions are generated from raw text using both rule-based mining patterns (Cheng et al., 2024b) and LLM-based question-answer generation (Jiang et al., 2024b). The raw documents are then combined with the generated instructions for continual training. Details of the training dataset construction are provided in Appendix E.1.

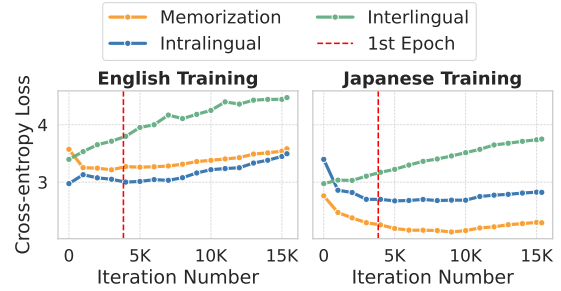
Training setup: We adopt llm-jp-3-13B (LLM-jp et al., 2024), a strong Japanese–English bilingual LLM, as the base model for pretraining, owing to its superior language ability in both languages, particularly Japanese. For each language, we cut off 0.5B tokens from the constructed corpus and train the data on llm-jp-3-13B for four epochs. Training hyperparameters are detailed in Appendix E.2.

3.2 Tracing Performance Dynamics

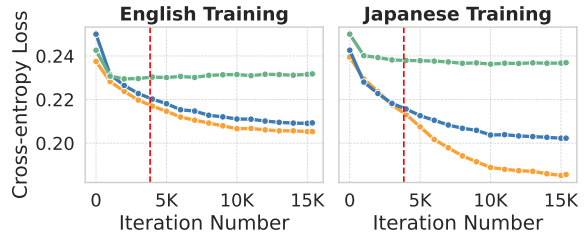
Figure 1 reports AdaXEval results for English and Japanese monolingual training. The results indicate that domain knowledge is gradually acquired during training in both languages.

Memorization evaluation: Accuracy increases from 36.1% to 59.6% (+23.5%) in English and from 45.7% to 75.9% (+30.2%) in Japanese. The higher post-training accuracy in Japanese partly reflects the stronger medical knowledge base of llm-jp-3-13B. However, since the factual instances used for evaluation differ between languages, direct cross-lingual comparison is not strictly fair.

Intralingual generalization: Both languages exhibit strong performance on the paraphrased datasets, with accuracy increasing from 32.6% to 54.7% (+22.1%) in English and from 34.9% to 53.8% (+18.9%) in Japanese. Notably, the improvement in English paraphrases parallels the



(a) Loss dynamics of correct query–answer sequences.



(b) Loss ratio dynamics for sequences with correct answers relative to all candidates.

Figure 2: Loss dynamics of datasets generated by AdaX-Eval during monolingual training.

memorization gain, whereas it is about 10% lower in Japanese, suggesting that the difficulty of intralingual generalization differs across languages.

Interlingual generalization: Figure 1 reveals that monolingual training results in limited cross-lingual knowledge transfer, yielding only 3.6% improvement in English-to-Japanese and 3.1% improvement in Japanese-to-English transfer.

3.3 Knowledge Acquisition via Loss Shielding

To elucidate the mechanisms of knowledge acquisition, this section analyzes the sequence loss of evaluation data to understand how knowledge is acquired during training. We analyze loss as it directly drives predictions on our multiple-choice datasets (see § 2.2) and reflects the model’s generation behavior, where sequences with lower loss are more likely to be generated.

(1) Training overfits to data, but the loss shielding drives knowledge memorization. We calculate the sequence loss of queries paired with correct answers in the evaluation dataset obtained by AdaXEval. Figure 2a shows the loss trajectory across training checkpoints for English and Japanese training. On the cloze prompt dataset, the loss decreases in early training but rises in later training, suggesting that training causes the model to overfit to the training corpus. However, Figure 1 reveals that memorization accuracy continues to

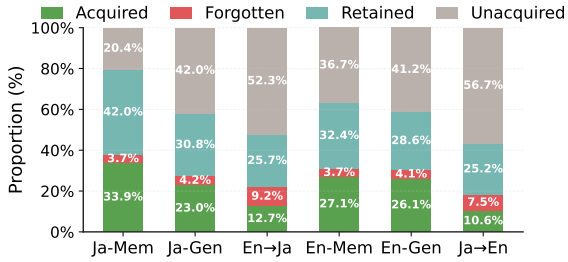


Figure 3: Instance state transitions before/after training.

improve until the 10K-th iteration. To investigate this divergence, we measure the ratio of the correct-sequence loss to the total loss across all four options. As shown in Figure 2b, this ratio mirrors the accuracy trend, suggesting that knowledge can still be memorized even under overfitting, since correct sequences are shielded from rapid loss growth, a phenomenon we term *loss shielding*. Figure 2b also explains the gap between cloze prompts and paraphrases, which the significant loss ratio gap in Japanese predicts, as shown in Figure 1.

(2) A trade-off exists between knowledge acquisition and forgetting. For each instance, we check its state transition before and after training by examining the loss. Figure 3 shows the proportions of instances retained, acquired, forgotten, or unacquired during training. Forgetting remains limited in monolingual evaluations, including both memorization and intralingual generalization. In contrast, cross-lingual transfer exhibits a notable increase in forgotten cases, offsetting gains from newly acquired knowledge. This suggests that while training in one language can introduce transferable knowledge, it also causes a decline in performance in other languages due to forgetting.

(3) Instance-level case studies show diverse knowledge acquisition patterns. We then examine the loss dynamics of instances acquired after training, analyzing all four options. We observe that the loss dynamics of correct answers follow distinct patterns: they either decrease steadily (Stable-Gain), increase while remaining lower than incorrect options (Loss-Shielding), or exhibit unstable behavior. Figure 4 illustrates three examples corresponding to the three loss change patterns.

4 From Training Data to Knowledge

This section investigates how knowledge of unseen queries is derived from the training data. Previous analyses have shown that the model tends to in-

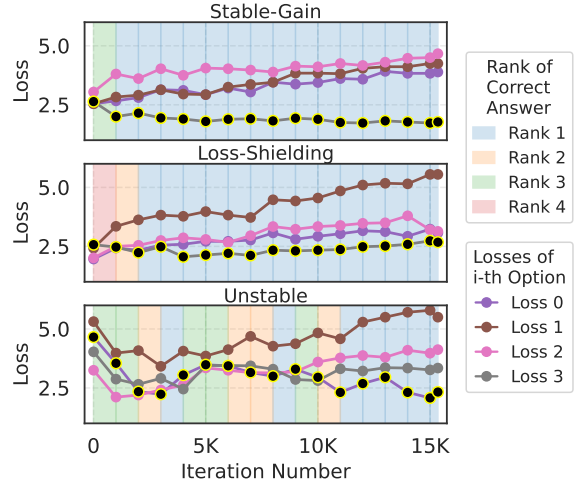


Figure 4: Acquired instances with different loss shapes. The line with bright circles indicates the correct answer.

crease the loss for sequences that express the same domain knowledge as the training data but in different linguistic forms. Although loss shielding allows the model to continue acquiring knowledge even as the loss for such sequences rises, this effect eventually fades once the loss grows too large. Understanding what sustains and what breaks this shielding effect is therefore essential for developing more robust training strategies. To this end, we introduce controlled perturbations into the training data by injecting noise under different rules, and track how their loss evolves during training. Specifically, we randomly sample 2,000 sequences from Japanese monolingual training data and apply perturbations at both the token and sequence levels.

4.1 Token-level Perturbation

Token sequences are perturbed after tokenization using the methods described below.

- **mask-X**: Replace X% tokens with <unk> token.
- **random-X**: Replace X% tokens with randomly sampled tokens from the tokenizer vocabulary.
- **delete-X**: Delete X% tokens.
- **reorder-X@Y**: Reorder tokens to achieve an edit distance equal to X% of the sequence length, with swaps restricted to a window of size Y.
- **monosyn-X**: Replace with Japanese synonyms.
- **mltsyn-X**: Replace with English synonyms.

Specifically, we set $X = 2^{1, \dots, 5}$ and $Y = 2^{0, \dots, 4}$. Synonym substitutions use the Japanese WordNet 2.0 (Bond and Kuribayashi, 2023) to retrieve Japanese or English synonyms. Additional implementation details are provided in Appendix D.

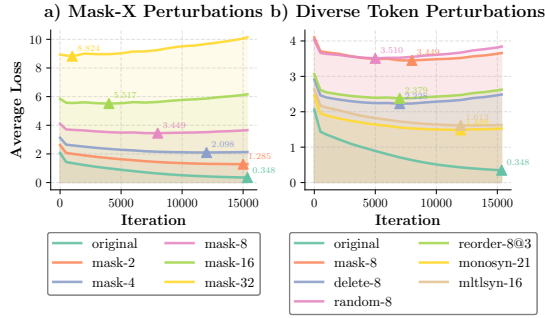


Figure 5: Loss dynamics of token-level perturbation.

(1) **More token edits cause larger loss harm.** Taking *mask-X* as an example, Figure 5(a) illustrates the loss dynamics across edited sequences compared to their original counterparts. Triangles (\blacktriangle) mark the *overfitting onset*, the point where minimum loss is attained before overfitting induces an upward trend. The results indicate that while the loss of original sequences decreases steadily, masking tokens introduce a substantial initial loss increase and accelerate the onset of overfitting. The analysis of other perturbation patterns leads to the same conclusion. See Appendix D for more results.

(2) **Loss sensitivity varies with vocabulary and structural perturbations.** Figure 5(b) illustrates the loss variations across sequences with 8% of tokens perturbed using diverse methods. Semantically aligned modifications (*monosyn* and *mltsyn*), exhibit the least impact on loss, followed by structural alterations (*reorder*, *delete*) that avoid introducing new vocabulary. Perturbations introducing irrelevant tokens (*mask*, *random*) inflict the greatest harm, significantly increasing the initial loss and accelerating the onset of overfitting.

4.2 Sentence-level Perturbation

We further perturb target sequences at the sentence level to simulate more realistic noise patterns, considering the following perturbation strategies.

- **partial-a**: Split each training document into four segments, then select the a -th segment.
- **syntax-X**: Rewrite $X\%$ sentences, modifying only syntax without changing vocabulary.
- **lexicon-X**: Rewrite $X\%$ sentences, modifying only vocabulary without changing syntax.
- **semantic-X**: Rewrite $X\%$ sentences, allowing both syntactic and lexical changes.
- **translation-X**: Translate $X\%$ sentences.

(1) **Losses exhibit stronger dependence on prior context.** Using partial sentences following the

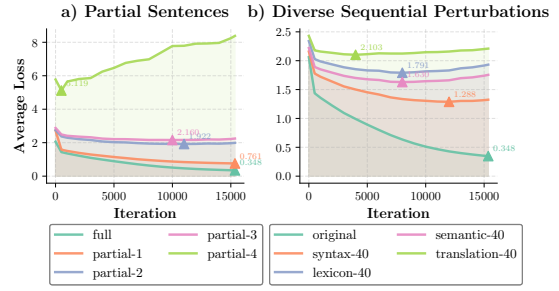


Figure 6: Loss dynamics of sequence-level perturbation.

partial-a strategy, Figure 6(a) shows that sentences appearing later in a training document incur higher loss when evaluated in isolation. This suggests that losses depend heavily on prior context, which constrains knowledge acquisition and highlights the need for a more balanced, context-independent learning paradigm.

(2) **Preserving vocabulary improves loss robustness.** Figure 6(b) presents the loss dynamics for all rewriting patterns applied to 40% of sentences, including both paraphrasing and translation. Among these patterns, syntax rewriting shows the most substantial loss shielding effect, followed by semantic rewriting. This is because semantic paraphrasing does not require extensive vocabulary replacement, resulting in fewer word substitutions compared to the lexicon paraphrases. These results highlight the necessity of incorporating lexicon-focused paraphrasing during training to improve models' ability to generalize knowledge across diverse test inputs, consistent with the observations in § 4.1.

5 Bridging Languages in Domain Adaptation

5.1 What Helps Cross-lingual Transfer?

We examine how cross-lingual transfer occurs by tracking the loss of target-language documents in models trained on source-language data. As the AdaXEval interlingual loss continues to increase (Figure 2a), tracing loss on evaluation data becomes unreliable. Instead, we examine the loss of the training documents across languages, serving as indirect evidence. Specifically, we sample 1,000 documents for each language and measure per-token loss on both datasets under two monolingual training settings. Since cross-lingual transfer converges rapidly within the first 1,000 iterations, we analyze the loss dynamics at a finer granularity during the first epoch. The results are shown in Figure 7. We observe that in-language loss consistently

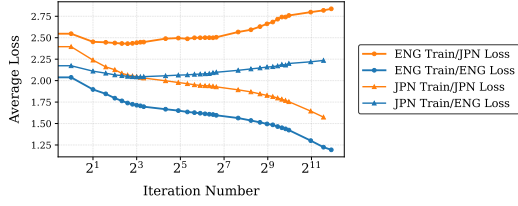


Figure 7: Losses on English/Japanese training data.

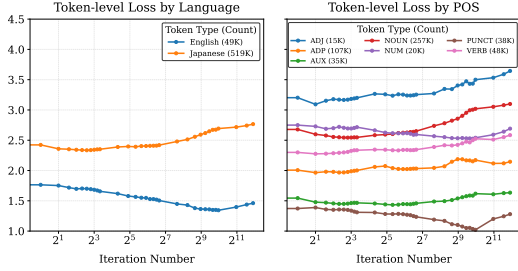


Figure 8: Japanese per-token loss by English training.

decreases, whereas cross-language loss decreases only during the initial iterations, reflecting the rapid but short-lived interlingual generalization.

To investigate what tokens contribute to the initial loss reduction, we further measure the per-token cross-entropy loss within the given sequences. Specifically, we examine the impact of two linguistic characteristics on loss changes: language and part-of-speech (POS). For language, tokens are classified by the language in which they occur, either Japanese or English; For POS, tokens are grouped by their syntactic category. Figure 8 shows token-level loss dynamics on Japanese documents during English training. Language-based token analysis indicates that English tokens embedded in Japanese documents benefit from English training. POS-based analysis supports the findings, as only numerical (NUM) and punctuation (PUNCT) tokens, both present in the English corpus, exhibit notable reductions in loss. This finding suggests that tokens occurring in the training corpus experience greater learning gains.

5.2 Cross-lingual Transfer Enhancement

In this section, we explore the factors in the training corpus that can drive improvements in cross-lingual transfer during domain adaptation. Due to the limited cross-lingual transfer ability in monolingual training, we switch to multilingual domain adaptation. Specifically, we investigate: what types of multilingual corpora can facilitate knowledge transfer during domain adaptation?

5.2.1 Multilingual Continual Training

To examine the factors that facilitate cross-lingual transfer, we construct a series of multilingual corpora and evaluate their effectiveness. To ensure a fair comparison, each training corpus is composed of two parts: a **knowledge injection corpus** \mathcal{C}_K and a **cross-lingual transfer enhancement corpus** \mathcal{C}_T . The \mathcal{C}_K contains the target knowledge expressed exclusively in the source language X , while \mathcal{C}_T , by contrast, does not provide any new knowledge in the target domain; instead, it serves only to establish linguistic connections between source language X and target language Y . Finally, we evaluate cross-lingual knowledge transfer by measuring the model’s ability to acquire knowledge in language Y , using the same evaluation metrics introduced in § 2.2. We fix the source language X as English and the target language Y as Japanese to facilitate analysis in this section. See Appendix F for the evaluation on reverse transfer direction.

Cross-lingual transfer enhancement corpora: We construct diverse corpora to enhance cross-lingual transfer using two primary strategies: *translation* and *romanization*. Both can build token-level connections between English and Japanese, where romanization requires significantly less effort for data collection. As baselines, we consider an empty cross-lingual corpus ($\mathcal{C}_T = \emptyset$) (**Monolingual**) and a strong domain-specific baseline using J-STAGE Japanese data, with documents related to the AdaX-Eval evaluation filtered out (**Medical-Japanese**).

To examine which translation data are most effective for domain adaptation, we prepare three types of bilingual corpora and use them to generate translation instructions as \mathcal{C}_T :

- **JParaCrawl (Balanced-Translation):** An English–Japanese web-crawled corpus covering diverse domains (Morishita et al., 2022).
- **ASPEC (Science-Translation):** A multilingual corpus containing academic paper abstracts across various scientific fields (Nakazawa et al., 2016). Documents in the medical and chemical domains are excluded to distinguish them from the target medical domain.
- **J-STAGE (Medical-Translation):** J-STAGE represents the closest domain to our target. We filter out all documents included in the AdaXEval evaluation datasets to avoid contamination.

To evaluate the romanization strategy, we construct a medical romanization dataset (**Medical-Roman**). We convert J-STAGE Japanese text to

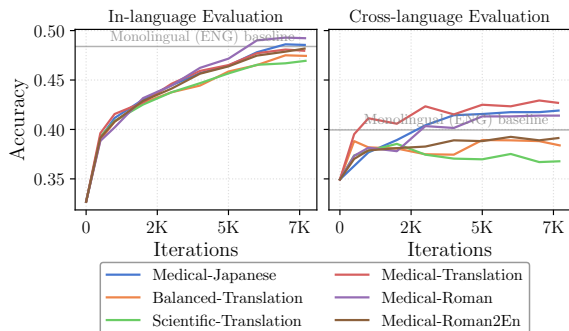


Figure 9: Enhancing cross-lingual transfer using diverse multilingual domain corpora.

romaji using `cutlet`³, an open-source tool for romanization. We then generate romanization instructions to link the Latin script of English with the Japanese script (kanji, hiragana, etc.). Finally, we create translation instructions between romanized Japanese and English based on J-STAGE (**Medical-Roman2En**), which serves as a comparison group. **Details of training:** Focusing on knowledge transfer from English to Japanese, we prepare two corpora, \mathcal{C}_K and \mathcal{C}_T , each containing 0.5 billion tokens, except for the **Monolingual** baseline. The English J-STAGE pretraining data serve as \mathcal{C}_K . For \mathcal{C}_T , we select seven candidate datasets as described in Sec. 5.2.1. We then combine the two 0.5B-tokens corpora into a single 1B-token corpus, shuffle it, and train `llm-jp-3-13B` with it for one epoch. See Appendix E for more training details.

5.2.2 Cross-lingual Transfer Evaluation

(1) Domain-specific corpus enhances cross-lingual transfer. Figure 9 presents the accuracy dynamics of multilingual domain adaptation, comparing two baselines with three translation datasets and two romanization datasets. Among these, only the corpus constructed from domain-specific data (**Medical-Japanese/Roman/Translation**) surpasses the **Monolingual** baseline, and only **Medical-Translation** achieves higher performance than both baselines. These results indicate that effective cross-lingual transfer of domain knowledge requires domain-specific signals; general cross-lingual enhancement methods fail to yield comparable improvements for domain knowledge. Figure 10 further illustrates the state transitions induced by cross-lingual transfer, showing that corpora yielding performance gains both increase acquired instances and reduce forgotten ones,

³<https://github.com/polm/cutlet>

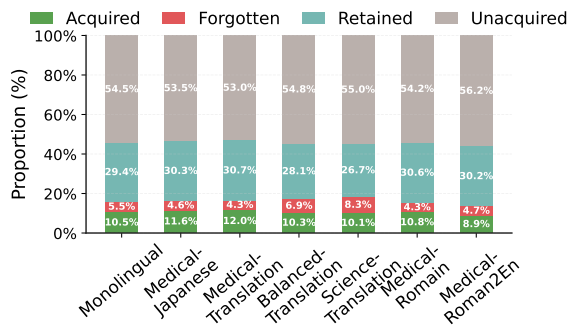


Figure 10: State transitions by cross-lingual transfer.

whereas the other datasets perform worse on both aspects. Finally, the stable performance across all recipes in the in-language evaluation indicates that using additional corpora unrelated to the target knowledge does not impair in-language knowledge acquisition.

(2) Achieving effective cross-lingual transfer of domain knowledge is challenging. Although both **Medical-Translation** yield improvements, the gains are still limited (around 3%) relative to the doubled training cost and additional dataset construction effort. This highlights the need for developing more efficient methods for cross-lingual domain knowledge transfer.

6 Conclusions

In this paper, we studied how LLMs acquire domain knowledge and transfer it across languages. We proposed **AdaXEval**, an adaptive evaluation pipeline that automatically generates datasets to assess domain knowledge across memorization, in-lingual generalization, and cross-lingual transfer. Using **AdaXEval**, we conduct a case study focusing on English-Japanese biomedical domain adaptation. We analyze the training dynamics of domain adaptation and find that knowledge acquisition is driven by **loss shielding**, wherein overfitting raises losses on irrelevant representations more rapidly than on relevant ones. Through **perturbation analysis**, we further reveal the sensitivity of LLMs to training data, offering practical guidance for developing more robust training paradigms. We also identified key factors that facilitate cross-lingual transfer through **controlled multilingual continual training**, showing that the presence of cross-lingual tokens in closely related domains is crucial.

As future work, we aim to develop a more training-robust and efficient pretraining paradigm to achieve domain knowledge acquisition.

7 Limitations

While this study provides novel insights into the mechanisms of bilingual domain adaptation, several limitations remain. First, our analysis is based on a single case study involving English–Japanese biomedical adaptation. Although this setting offers a controlled environment to examine multilingual knowledge acquisition, the generalizability of the findings to other domains, language pairs, and model architectures remains to be validated. Future research should therefore extend the investigation to a broader range of large language models, including those trained under different configurations, as well as to diverse domain and linguistic contexts.

Second, our evaluation primarily relies on loss-based quantitative metrics to characterize the dynamics of knowledge acquisition and forgetting. While these measures provide a consistent and interpretable framework for tracking learning behavior, they do not fully capture how such knowledge is manifested in model outputs. A complementary qualitative and text-level evaluation, examining generated outputs and factual correctness, will be necessary to bridge the gap between internal training dynamics and externally observable performance.

References

Francis Bond and Takayuki Kuribayashi. 2023. *The Japanese Wordnet 2.0*. In *Proceedings of the 12th Global Wordnet Conference*, pages 179–186, University of the Basque Country, Donostia - San Sebastian, Basque Country. Global Wordnet Association.

Hoyeon Chang, Jinho Park, Seonghyeon Ye, Sohee Yang, Youngkyung Seo, Du-Seong Chang, and Minjoon Seo. 2024. *How do large language models acquire factual knowledge during pretraining?* *Preprint*, arXiv:2406.11813.

Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. 2023. *Journey to the center of the knowledge neurons: Discoveries of language-independent knowledge neurons and degenerate knowledge neurons*. *Preprint*, arXiv:2308.13198.

Daixuan Cheng, Yuxian Gu, Shaohan Huang, Junyu Bi, Minlie Huang, and Furu Wei. 2024a. *Instruction pre-training: Language models are supervised multitask learners*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2529–2550, Miami, Florida, USA. Association for Computational Linguistics.

Daixuan Cheng, Shaohan Huang, and Furu Wei. 2024b. *Adapting large language models via reading compre-*

hension. In *The Twelfth International Conference on Learning Representations*.

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. *Knowledge neurons in pretrained transformers*. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8493–8502, Dublin, Ireland. Association for Computational Linguistics.

DeepSeek-AI. 2025. *Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning*. *Preprint*, arXiv:2501.12948.

Yihao Fang, Xianzhi Li, Stephen Thomas, and Xiaodan Zhu. 2023. *ChatGPT as data augmentation for compositional generalization: A case study in open intent detection*. In *Proceedings of the Fifth Workshop on Financial Technology and Natural Language Processing and the Second Multimodal AI For Financial Forecasting*, pages 13–33, Macao. -.

Changjiang Gao, Hongda Hu, Peng Hu, Jiajun Chen, Jixing Li, and Shujian Huang. 2024a. *Multilingual pre-training and instruction tuning improve cross-lingual knowledge alignment, but only shallowly*. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6101–6117, Mexico City, Mexico. Association for Computational Linguistics.

Changjiang Gao, Hongda Hu, Peng Hu, Jiajun Chen, Jixing Li, and Shujian Huang. 2024b. *Multilingual pre-training and instruction tuning improve cross-lingual knowledge alignment, but only shallowly*. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6101–6117, Mexico City, Mexico. Association for Computational Linguistics.

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, and 5 others. 2023. *A framework for few-shot language model evaluation*.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, and Others. 2024. *The llama 3 herd of models*. *Preprint*, arXiv:2407.21783.

Evan Hernandez, Belinda Z Li, and Jacob Andreas. 2023. *Measuring and manipulating knowledge representations in language models*. *arXiv preprint arXiv:2304.00740*.

Seongtae Hong, Seungyoon Lee, Hyeonseok Moon, and Heuseok Lim. 2025. *MIGRATE: Cross-lingual adaptation of domain-specific LLMs through code-switching and embedding transfer*. In *Proceedings of*

752	<i>the 31st International Conference on Computational Linguistics</i> , pages 9184–9193, Abu Dhabi, UAE. Association for Computational Linguistics.	
753		
754		
755	Ryosuke Ishigami. 2025. Deepseek-r1-distill-qwen-32b-japanese .	
756		
757	Jaavid J, Raj Dabre, Aswanth M, Jay Gala, Thanmay Jayakumar, Ratish Puduppully, and Anoop Kunchukuttan. 2024. RomanSetu: Efficiently unlocking multilingual capabilities of large language models via Romanization . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15593–15615, Bangkok, Thailand. Association for Computational Linguistics.	
758		
759		
760		
761		
762		
763		
764		
765		
766	Joel Jang, Seonghyeon Ye, Changho Lee, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, and Minjoon Seo. 2022a. TemporalWiki: A lifelong benchmark for training and evaluating ever-evolving language models . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 6237–6250, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
767		
768		
769		
770		
771		
772		
773		
774		
775	Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, Stanley Jungkyu Choi, and Minjoon Seo. 2022b. Towards continual knowledge learning of language models . <i>Preprint</i> , arXiv:2110.03215.	
776		
777		
778		
779		
780	Junfeng Jiang, Fei Cheng, and Akiko Aizawa. 2024a. Improving referring ability for biomedical language models . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 6444–6457, Miami, Florida, USA. Association for Computational Linguistics.	
781		
782		
783		
784		
785		
786	Junfeng Jiang, Jiahao Huang, and Akiko Aizawa. 2025. JMedBench: A benchmark for evaluating Japanese biomedical large language models . In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 5918–5935, Abu Dhabi, UAE. Association for Computational Linguistics.	
787		
788		
789		
790		
791		
792	Ting Jiang, Shaohan Huang, Shengyue Luo, Zihan Zhang, Haizhen Huang, Furu Wei, Weiwei Deng, Feng Sun, Qi Zhang, Deqing Wang, and Fuzhen Zhuang. 2024b. Improving domain adaptation through extended-text reading comprehension . <i>Preprint</i> , arXiv:2401.07284.	
793		
794		
795		
796		
797		
798	Mathieu Lai-king and Patrick Paroubek. 2024. Pre-training data selection for biomedical domain adaptation using journal impact metrics . <i>Preprint</i> , arXiv:2409.02725.	
799		
800		
801		
802	Tian Liang, Xing Wang, Mingming Yang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2024. Addressing entity translation problem via translation difficulty and context diversity . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 11628–11638, Bangkok, Thailand. Association for Computational Linguistics.	
803		
804		
805		
806		
807		
808		
	Peiqin Lin, Andre Martins, and Hinrich Schuetze. 2025. A recipe of parallel corpora exploitation for multilingual large language models . In <i>Findings of the Association for Computational Linguistics: NAACL 2025</i> , pages 4038–4050, Albuquerque, New Mexico. Association for Computational Linguistics.	809
		810
		811
		812
		813
		814
	Yihong Liu, Mingyang Wang, Amir Hossein Kargaran, Felicia Körner, Ercong Nie, Barbara Plank, François Yvon, and Hinrich Schütze. 2025. Tracing multilingual factual knowledge acquisition in pretraining . <i>Preprint</i> , arXiv:2505.14824.	815
		816
		817
		818
		819
	LLM-jp, :, Akiko Aizawa, Eiji Aramaki, Bowen Chen, Fei Cheng, Hiroyuki Deguchi, Rintaro Enomoto, Kazuki Fujii, Kensuke Fukumoto, Takuya Fukushima, Namgi Han, and 1 others. 2024. Llm-jp: A cross-organizational project for the research and development of fully open japanese llms . <i>Preprint</i> , arXiv:2407.03963.	820
		821
		822
		823
		824
		825
		826
	Qwen Team. 2025. Qwen3 technical report . <i>Preprint</i> , arXiv:2505.09388.	827
		828
	Soumen Kumar Mondal, Sayambhu Sen, Abhishek Singhanian, and Preethi Jyothi. 2025. Language-specific neurons do not facilitate cross-lingual transfer . In <i>The Sixth Workshop on Insights from Negative Results in NLP</i> , pages 46–62, Albuquerque, New Mexico. Association for Computational Linguistics.	829
		830
		831
		832
		833
		834
	Makoto Morishita, Katsuki Chousa, Jun Suzuki, and Masaaki Nagata. 2022. JParaCrawl v3.0: A large-scale English-Japanese parallel corpus . In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 6704–6710, Marseille, France. European Language Resources Association.	835
		836
		837
		838
		839
		840
	Toshiaki Nakazawa, Manabu Yaguchi, Kiyotaka Uchi-moto, Masao Utiyama, Eiichiro Sumita, Sadao Kurohashi, and Hitoshi Isahara. 2016. ASPEC: Asian scientific paper excerpt corpus . In <i>Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)</i> , pages 2204–2208, Portorož, Slovenia. European Language Resources Association (ELRA).	841
		842
		843
		844
		845
		846
		847
		848
	Mark Neumann, Daniel King, Iz Beltagy, and Waleed Ammar. 2019. ScispaCy: Fast and Robust Models for Biomedical Natural Language Processing . In <i>Proceedings of the 18th BioNLP Workshop and Shared Task</i> , pages 319–327, Florence, Italy. Association for Computational Linguistics.	849
		850
		851
		852
		853
		854
	Jingcheng Niu, Andrew Liu, Zining Zhu, and Gerald Penn. 2024. What does the knowledge neuron thesis have to do with knowledge? In <i>The Twelfth International Conference on Learning Representations</i> .	855
		856
		857
		858
	Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. 2019. Language models as knowledge bases? <i>arXiv preprint arXiv:1909.01066</i> .	859
		860
		861
		862
	Sukannya Purkayastha, Sebastian Ruder, Jonas Pfeiffer, Iryna Gurevych, and Ivan Vulić. 2023.	863
		864

865	Romanization-based large-scale adaptation of multilingual language models.	Zhihan Zhang, Dong-Ho Lee, Yuwei Fang, Wenhao Yu, Mengzhao Jia, Meng Jiang, and Francesco Barbieri. 2024. PLUG: Leveraging pivot language in cross-lingual instruction tuning.	922
866			923
867			924
868			925
869			926
870	Uri Shaham, Jonathan Herzig, Roei Aharoni, Idan Szpektor, Reut Tsarfaty, and Matan Eyal. 2024. Multilingual instruction tuning with just a pinch of multilinguality.	Jun Zhao, Zhihao Zhang, Luhui Gao, Qi Zhang, Tao Gui, and Xuanjing Huang. 2024a. Llama beyond english: An empirical study on language capability transfer.	927
871			928
872			929
873			930
874			931
875			932
876			933
877	Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Scharli, Aakanksha Chowdhery, Philip Mansfield, Blaise Aguerre y Arcas, Dale Webster, and 11 others. 2022. Large language models encode clinical knowledge.	Xin Zhao, Zehui Jiang, and Naoki Yoshinaga. 2025. Neuron empirical gradient: Discovering and quantifying neurons’ global linear controllability.	934
878			935
879			936
880			937
881			938
882			939
883			940
884			941
885			942
886	Social Computing Lab. 2023. Mednern-cr-ja (revision 13dbcb6).	Xin Zhao, Naoki Yoshinaga, and Daisuke Oba. 2024b. Tracing the roots of facts in multilingual language models: Independent, shared, and transferred knowledge.	943
887			944
888			945
889	Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. 2022. Finding skill neurons in pre-trained transformer-based language models.		946
890			947
891			948
892			949
893			950
894			951
895	Orgest Xhelili, Yihong Liu, and Hinrich Schuetze. 2024. Breaking the script barrier in multilingual pre-trained language models with transliteration-based post-training alignment.	Xin Zhao, Naoki Yoshinaga, and Daisuke Oba. 2024c. What matters in memorizing and recalling facts? multifaceted benchmarks for knowledge probing in language models.	952
896			953
897			954
898			955
899			956
900			957
901			958
902			959
903			960
904	Qianqian Xie, Weiguang Han, Zhengyu Chen, Ruoyu Xiang, Xiao Zhang, Yueru He, Mengxi Xiao, Dong Li, Yongfu Dai, Duanyu Feng, Yijing Xu, Haoqiang Kang, Ziyang Kuang, Chenhan Yuan, Kailai Yang, Zheheng Luo, Tianlin Zhang, Zhiwei Liu, Guojun Xiong, and 15 others. 2024a. Finben: A holistic financial benchmark for large language models.	Nicolas Zucchet, Jörg Bornschein, Stephanie Chan, Andrew Lampinen, Razvan Pascanu, and Soham De. 2025. How do language models learn facts? dynamics, curricula and hallucinations.	961
905			962
906			963
907			964
908			965
909			966
910	Yong Xie, Karan Aggarwal, and Aitzaz Ahmad. 2024b. Efficient continual pre-training for building domain specific large language models.		967
911			968
912			969
913			970
914			971
915			972
916	Ikuya Yamada and Ryokan Ri. 2024. LEIA: Facilitating cross-lingual knowledge transfer in language models with entity-based data augmentation.		973
917			974
918			975
919			
920			
921			

976 not necessarily reflect domain-specific knowledge
977 that users care about.

978 A.2 Mechanisms of Knowledge Acquisition

979 LLMs have been widely studied for their capacity
980 to store and retrieve factual knowledge (Petroni
981 et al., 2019; Hernandez et al., 2023; Chang et al.,
982 2024; Zhao et al., 2024c), spurring interest in un-
983 derstanding how such knowledge is encoded and
984 accessed (Dai et al., 2022; Wang et al., 2022; Niu
985 et al., 2024; Zhao et al., 2025). Recent work has
986 investigated knowledge acquisition dynamics by
987 analyzing intermediate checkpoints during train-
988 ing. Chang et al. (2024) shows that LLMs accu-
989 mulate factual knowledge through repeated expo-
990 sures, gradually increasing recall likelihood with
991 each encounter. Key factors influencing acquisi-
992 tion include fact frequency, model scale, and batch
993 size (Liu et al., 2025; Zhao et al., 2024b).

994 Meanwhile, cross-lingual transfer ability has
995 also garnered considerable attention for enabling
996 efficient knowledge acquisition across languages.
997 However, transferring factual knowledge across
998 languages remains notably challenging (Liu et al.,
999 2025; Zhao et al., 2024b). Facts expressed in differ-
1000 ent languages may be stored as distinct representa-
1001 tions (Chen et al., 2023; Zhao et al., 2024b; Mondal
1002 et al., 2025), and only certain relational types trans-
1003 fer effectively (Liu et al., 2025; Zhao et al., 2024b).
1004 Despite these advances, existing studies focus on
1005 predefined relational facts (Zhao et al., 2024b; Liu
1006 et al., 2025), which inadequately represent com-
1007 plex knowledge in specialized domains. Moreover,
1008 the connection between training data and acquired
1009 knowledge remains underexplored, which is criti-
1010 cal for designing optimal training recipes.

1011 A.3 Cross-lingual Knowledge Transfer

1012 Cross-lingual transfer in domain adaptation has
1013 been studied through various strategies that aim
1014 to bridge language gaps in knowledge transfer.
1015 Translation data are widely used in training LLMs
1016 to enhance cross-lingual knowledge transfer (Lin
1017 et al., 2025; Zhao et al., 2024a; Gao et al., 2024b),
1018 supporting applications such as machine transla-
1019 tion (Zhao et al., 2024a) and instruction follow-
1020 ing (Shaham et al., 2024; Zhang et al., 2024). How-
1021 ever, since in-domain translation data are often
1022 scarce, it remains unclear whether translation can
1023 effectively transfer complex and sparse domain
1024 knowledge across languages. Furthermore, the ro-
1025 manization strategy reduces script barriers by con-

verting text into romanized forms, thereby align- 1026
ing linguistic representations with English (J et al., 1027
2024; Purkayastha et al., 2023; Xhelili et al., 2024). 1028
Finally, code-switching has been shown to be effec- 1029
tive in aligning token semantics across languages 1030
by interleaving tokens from multiple languages 1031
within the same context (Hong et al., 2025; Ya- 1032
mada and Ri, 2024). However, it is still uncertain 1033
whether these alignment signals are sufficient for 1034
transferring complex domain knowledge across lan- 1035
guages. 1036

B AdaXEval Implementation Details 1037

B.1 Generation Details 1038

In this section, we provide details on the process 1039
used to generate the AdaXEval dataset from J- 1040
STAGE documents. 1041

(1) **Factual sentence filtering:** We use the open- 1042
source NLP tool HanLP⁴ for Japanese sentence seg- 1043
mentation, and scispaCy⁵ (Neumann et al., 2019) 1044
, which provides a full spaCy pipeline for scien- 1045
tific and biomedical texts, for English sentence seg- 1046
mentation. Subsequently, we perform biomedical 1047
named entity recognition (NER) on each sentence 1048
and filter out sentences containing fewer than two 1049
named entities. Specifically, for Japanese medical 1050
documents, we employ MedNERN-CR-JA⁶ (So- 1051
cial Computing Lab, 2023) , a model specialized 1052
for NER in the Japanese medical domain. For En- 1053
glish texts, we use the “en_ner_bionlp13cg_md” 1054
model provided by scispaCy for biomedical entity 1055
extraction. 1056

(2) **Domain triple extraction:** In the next step, 1057
we use a multi-LLM agent to first judge whether 1058
the given sentence contains biomedical facts and 1059
extract them if it does. We carefully design CoT 1060
instruction with few-shot examples to instruct each 1061
LLM in the agent to generate a structural output, 1062
containing three fields: 1063

- factuality: answer yes or no. If yes, the 1064
model should output the triple; Otherwise, it 1065
outputs None to the triple field. 1066
- triple: A nested JSON object with three fields: 1067
subject, relation, and object representing the ex- 1068
tracted fact. Please note that if factuality is false, 1069
make sure this field is null 1070

⁴<https://github.com/hankcs/HanLP>

⁵<https://github.com/allenai/scispaCy>

⁶<https://github.com/sociocom/MedNERN-CR-JA>

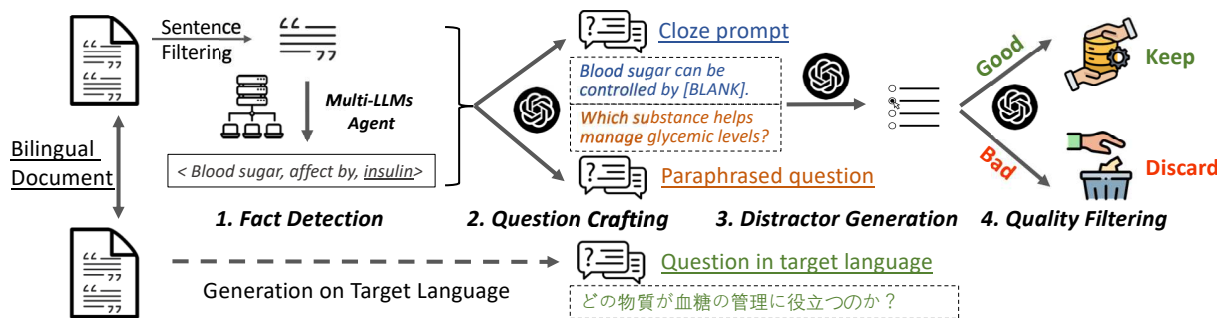


Figure 11: Overview of AdaXEval, a pipeline to adaptively generate domain knowledge evaluation datasets.

- reason: A brief explanation for why the sentence was or wasn't considered factual, referring to the criteria provided. This is included to improve the model's reasoning ability.

The models' confidence in judging the factuality is measured by the probability of yes token out by the factuality field.

For our experiments, we utilize three strong open-source LLMs for both English and Japanese, opting for open-source models to avoid the high computational cost of commercial APIs. For **English biomedical triple extraction**, we employ Qwen-32B (Qwen Team, 2025),⁷ DeepSeek-R1-Distill-Llama-70B (DeepSeek-AI, 2025),⁸ and Llama-3.3-70B-Instruct (Grattafiori et al., 2024).⁹ For **Japanese biomedical triple extraction**, we use Qwen-32B, llm-jp-3.1-8x13b-instruct4,¹⁰ and Llama-3.3-Swallow-70B-Instruct-v0.4.¹¹ For each sentence, we aggregate the confidence scores from the three models and retain sentences with at least two models predicting a yes label. We then apply a heuristic method to select the final triple from the three candidates.

(3) Generation of queries and distractors: We finally use the extracted triples and their corresponding context sentences as input to instruct the strong close-source LLM, GPT-4.1, for generating queries and distractors. We carefully design the prompts for both English and Japanese. The final generation is recorded and included in the final AdaXEval evaluation dataset.

(4) Quality filtering: Finally, we use LLM to con-

⁷<https://huggingface.co/Qwen/Qwen3-32B>

⁸<https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-70B>

⁹<https://huggingface.co/meta-llama/Llama-3.3-70B-Instruct>

¹⁰<https://huggingface.co/llm-jp/llm-jp-3.1-8x13b-instruct4>

¹¹<https://huggingface.co/tokyotech-llm/Llama-3.3-Swallow-70B-Instruct-v0.4>

duct a two-steps quality filtering. First, we introduce a rigorous annotation instruction for assessing the quality of automatically generated biomedical knowledge questions. Given a sentence, its knowledge triple, a fill-in-the-blank query, and a question-style paraphrase, we evaluate three aspects:

- The fidelity and clarity of the cloze prompt,
- The semantic equivalence and self-containment of the paraphrased question, and
- The correctness and plausibility of the answer and distractors.

Then, we craft the annotation instruction to let LLM decide whether the document in the target languages contains the knowledge matching the created instances in the source languages. The annotation is to confirm the quality of the AdaXEval dataset in providing interlingual generalization evaluation.

B.2 Evaluation Metrics

We follow (Gao et al., 2023) to compute the average cross-entropy loss over the target tokens of possible answers and select the one that has the highest generation possibility as the final answer. Specifically, for loss calculation of cloze queries, we use tokens before the [BLANK] as context and compute loss on the following tokens. For paraphrases, we treat the question as context and measure only the loss of answer tokens. We use prediction accuracy as the metric for knowledge acquisition. **Formulation:**

Let the model be $p_\theta(\cdot | \cdot)$. Each dataset \mathcal{D} contains pairs (q, a) , where q is either a *cloze prompt* or a *paraphrase question*, and $a = (a_1, \dots, a_m)$ is the tokenized answer sequence (e.g., "insulin"). For cloze prompts, the prompt contains a special token [BLANK], and for paraphrases, the question is a natural question. We denote by c the context tokens and by $s = (s_1, \dots, s_n)$ the evaluation sequence whose loss we measure.

Cloze queries: For a cloze prompt such as:
 "[BLANK] can be used to control blood
 sugar level."
 the evaluation sequence s is the full completion
 after the [BLANK], i.e.,

$$s = (a_1, \dots, a_m, r_1, \dots, r_k),$$

where a_1, \dots, a_m are answer tokens ("insulin"),
 and r_1, \dots, r_k are the remainder tokens (" can be
 used to control blood sugar level"). The average
 cross-entropy loss is defined as

$$\mathcal{L}_{\text{cloze}}(q, a) = -\frac{1}{n} \sum_{t=1}^n \log p_{\theta}(s_t | c, s_{<t}).$$

Paraphrase queries: For a paraphrase question
 such as
 "Which substance helps manage glycemic
 levels?"

we take the entire question tokens as context c , and
 the evaluation sequence is only the answer tokens:

$$s = (a_1, \dots, a_m).$$

The loss is computed as

$$\mathcal{L}_{\text{para}}(q, a) = -\frac{1}{m} \sum_{t=1}^m \log p_{\theta}(a_t | c, a_{<t}).$$

Prediction and accuracy: For multiple-choice
 answers $\mathcal{A} = \{a^{(1)}, \dots, a^{(K)}\}$, we select the candi-
 date with the lowest loss (equivalently, highest
 likelihood):

$$\hat{a} = \arg \min_{a \in \mathcal{A}} \mathcal{L}(q, a).$$

Finally, the knowledge acquisition metric is accu-
 racy:

$$\text{Accuracy} = \frac{1}{|\mathcal{D}|} \sum_{(q,a) \in \mathcal{D}} \mathbf{1}[\hat{a} = a].$$

C AdaXEval Dataset Details

We randomly sampled 10,000 parallel documents
 to generate the evaluation dataset. The number of
 instances after each step in AdaXEval generation
 is shown in Table 1.

Table 1: Dataset statistics at each step of AdaXEval.

Step	English	Japanese
Sampled abstracts	10,000	10,000
Splitted sentences	81,770	71,661
Sentences after entity filtering (≥ 2 entities)	45,390	40,762
Triple extraction	4840	3926
Cloze queries generation	4840	3926
Paraphrases generation	4840	3926
After Quality Filtering	3236	2553

C.1 Human Evaluation

To assess the quality of our generated datasets, we
 conduct a comprehensive human evaluation across
 four key components of the knowledge extraction
 and question generation pipeline. The annotation
 is conducted by the first author, who has a lan-
 guage background of both Japanese and English.
 For the annotation that requires specific domain
 knowledge, the author uses advanced LLMs, such
 as ChatGPT or Claude, as an assistant for anno-
 tation. See the full human evaluation guideline
 documented in Appendix G.

(1) Cloze prompt evaluation checks the faithful-
 ness of the generated prompt to the original sen-
 tence structure.

(2) Paraphrases evaluation is conducted on four
 dimensions: fluency and grammaticality, linguistic
 diversity in reformulation, factual correctness to the
 original sentence in the source language, and inter-
 lingual factual correctness using the corresponding
 documents in the target language.

(3) Distractor quality is measured through plausi-
 bility within the domain and apparent incorrectness
 relative to the original context.

Each metric employs structured scoring rubrics
 with scales ranging from 0-2 or 0-3, enabling sys-
 tematic assessment of dataset quality across mul-
 tiple linguistic and semantic dimensions. We ran-
 domly sample 50 instances for each language and
 conduct human evaluation following the guidance
 above. As shown in Table 2, our evaluation results
 indicate that AdaXEval is capable of generating
 high-quality evaluation data, meeting the require-
 ments for assessing knowledge memorization as
 well as intralingual and interlingual generalization
 evaluation.

Table 2: Human evaluation results for knowledge acquisition datasets generated from the biomedical J-STAGE corpus.

Evaluation Metric	Japanese	English
cloze prompt (Faithfulness)	2.84/3	2.89/3
Paraphrase (Fluency)	2.94/3	2.96/3
Paraphrase (Diversity)	2.48/3	2.62/3
Paraphrase (Factuality)	2.86/3	2.86/3
Paraphrase (Inter-Factuality)	1.68/2	1.76/2
Distractor (Plausibility)	2.35/3	2.16/3
Distractor (Incorrectness)	2.89/3	2.97/3

C.2 Examples

We randomly sample 10 examples from AdaXEval for both English and Japanese and display them in Table 3 (Japanese) and Table 4 (English).

D Sequence Perturbation Analysis

In this section, we introduce the detailed settings for sequence perturbation experiments and report the additional results.

D.1 Details of Perturbation

For *monosym@X* and *monosym@X* that require collecting synonyms from WordNet 2.0, we only conduct Japanese-to-English replacement. Specifically, we first tokenize the Japanese sequence by *sudachipy*¹², a Japanese morphological analyzer, and get the POS tags. Then we filter out tokens with stop words and the POS tags that are not “普通名詞”, “固有名詞”, “サ変接続”, “形容動詞語幹”, “動詞一般” to avoid introducing noisy words. Furthermore, the paraphrasing and translation are done by requesting GPT-4.1.

D.2 Perturbation Results

We show all the token-level perturbation results in Figure 13 and sentence-level perturbation result in Figure 14.

E Domain Adaptation Training Details

E.1 Training Data Generation

In this study, to address the scarcity of bilingual domain corpora and enhance domain understanding, we employ two data augmentation strategies: regex-based pattern mining and LLM-based QA generation. Both approaches yield instruction-like

sequences, which we mix with raw corpora. Following prior work (Cheng et al., 2024a), we adopt an instruction-pretraining strategy for continual domain adaptation.

(1) Regex-based pattern mining: (Jiang et al., 2024b) verified that by transforming raw corpora into reading comprehension texts, continual training can consistently enhance performance across various tasks in different domains. We adopt a similar strategy by analyzing the training corpora and mining regex patterns to automatically create instruction-style data. Furthermore, to increase data diversity, we prepare ten instruction templates for each type of reading comprehension text. For each document, however, we sample only one template per type.

Specifically, each document in J-STAGE contains multiple metadata fields, including:

- title: the title of the paper
- abstract: the paper abstract
- keywords: pre-defined keywords of the paper
- fields: research categories of the paper

Based on this information, we construct ten types of reading comprehension instructions as follows:

- **Summarization:** Summarize the context into one concise sentence, taking the abstract as input and the title as output.
- **Keyword Extraction:** Extract the keywords from the abstract, using the keywords field as the gold reference.
- **Field Identification:** Identify the research field(s) of the paper, taking the abstract as input and the fields metadata as the expected output.
- **Translation:** Translate between English and Japanese, using bilingual metadata or parallel text segments as input-output pairs.
- **Text Completion:** Complete an incomplete abstract or title given the partial text, where the remainder of the text serves as the reference output.
- **Conclusion Derivation:** Derive the study’s conclusion from its context, with the conclusion section as supervision.
- **Background Derivation:** Infer the background or motivation of the study from the provided abstract or introduction sentences.

¹²<https://github.com/WorksApplications/SudachiPy>

- **Diagnosis:** Given a description of symptoms (extracted from biomedical corpora), predict the corresponding diagnosis, using annotated datasets where available.
- **Reordering:** Reorder shuffled sentences into their natural sequence, ensuring coherence with the original abstract or section structure.
- **Goal–Method–Result–Conclusion (GMRC):** Derive one missing component (*e.g.*, goal, method, result, or conclusion) based on the other three, enabling comprehension of scientific discourse structures.

(2) **LLM-base QA generation:** To enhance the data diversity, we further generate five question-answer pairs for each document. Specifically, we use DeepSeek-R1-Distill-Llama-70B for English QA-pair generation and use DeepSeek-R1-Distill-Qwen-JP-32B (Ishigami, 2025)¹³ for Japanese QA generation. Noted that not all documents can successfully generate five QA pairs.

E.2 continual training Settings

We conduct continual training using **Megatron-LM** on the `llm-jp-3-13B` model. The training setup follows a distributed configuration with 4 compute nodes, each equipped with 8 A100 GPUs. We apply a tensor parallel size of 2 and a pipeline parallel size of 4, enabling efficient large-scale training with a sequence length of 4096. The optimizer is configured with a learning rate of 2×10^{-5} , weight decay of 0.1, and gradient clipping of 1.0, with a minimum learning rate of 2×10^{-6} . We adopt a micro-batch size of 1 and a global batch size of 32 to stabilize training. Under this configuration, training one epoch on 0.5B tokens requires approximately 7 hours, demonstrating the computational feasibility of continual training while maintaining efficiency on large-scale biomedical and cross-lingual corpora.

F Transfer From Japanese to English

Here, we report the results of the Japanese-to-English cross-lingual transfer evaluation using different recipes, as a supplement to the analysis in § 5.2.2. Specifically, in this setting, \mathcal{C}_K is composed of the Japanese monolingual training corpora, while we vary the data source of \mathcal{C}_T to investigate efficient cross-lingual transfer. Evaluation results are shown in Figure 12.

¹³<https://huggingface.co/cyberagent/DeepSeek-R1-Distill-Qwen-32B-Japanese>

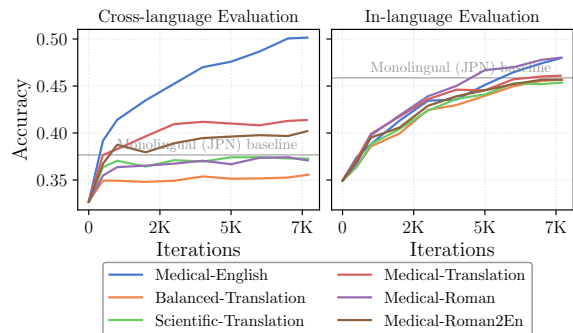


Figure 12: Japanese-to-English transfer evaluation with diverse strategies.

The evaluation results shown in Figure 12 exhibit trends that differ from those in Figure 9, as the strong baseline using multilingual corpora (a mixture of Japanese and English monolingual training corpora) largely outperforms other translation corpora. This suggests that generalization within English knowledge is an easier task than transferring knowledge from Japanese to English. Nevertheless, we still observe the strongest cross-lingual transfer in the closely related domain compared to the other three domains.

G AdaXEval Human Evaluation Guideline

G.1 Cloze prompt Quality Evaluation

Objective:

Evaluate whether fill-in-the-blank prompts generated from biomedical academic sentences and subject–relation–object triples are clear, faithful to the original sentence, and free from factual distortion. Each prompt removes the object from the triple and replaces it with a blank.

Input: You will be shown the following for each item:

- **Original Sentence:** The full academic sentence from which the triple is extracted.
- **Triple:** A <subject, relation, object> triple derived from that sentence.
- **Generated Prompt:** A fill-in-the-blank sentence where the object is replaced with [BLANK].

Faithfulness Criteria

When the [BLANK] is replaced with the object, does the prompt preserve the structure and meaning of the original sentence, including key contextual information?

- **What it checks:** Structural and contextual simi-

Score	Description
3	The prompt closely mirrors the original sentence's structure and meaning. Minor surface-level changes (<i>e.g.</i> , auxiliary verbs, punctuation, or sentence breaks) are acceptable.
2	The core meaning is preserved, but there are moderate changes in wording, omission that does not change the semantic meaning, or noticeable rephrasing.
1	The prompt differs significantly in structure or phrasing, or some key information is missed
0	The prompt is not clearly based on the original sentence or appears unrelated in meaning or form.

1370 larity between the prompt and the original sentence
1371

- 1372 • **Focuses on:** Wording, phrasing, sentence structure, presence of supporting context
1373

1374 G.2 Paraphrased Question Quality 1375 Evaluation

1376 **Objective:** Evaluate whether **question-style prompts**, automatically generated from fill-in-the-blank biomedical prompts, are:
1377
1378

- 1379 • Semantically faithful to the original prompt (*i.e.*, same question being asked)
1380
- 1381 • Grammatically correct and fluent
- 1382 • Natural as questions a human would realistically ask
1383

1384 Input:

- 1385 • Original Fill-in-the-Blank Prompt (*e.g.*, “EGFR is highly expressed in [BLANK].”)
1386
- 1387 • Paraphrased Question-style Prompt (*e.g.*, “In which condition is EGFR highly expressed?”)
1388

1389 Example:

- 1390 • Original Sentence: “EGFR is highly expressed in non-small cell lung carcinoma.”
1391
- 1392 • Triple: (*EGFR, is highly expressed in, non-small cell lung carcinoma*)
1393
- 1394 • Prompt: “EGFR is highly expressed in [BLANK].”
1395

1396 Evaluation Criteria

1397 **1) Fluency and Grammaticality:** Is the question grammatically correct, fluent, and natural-sounding in English?
1398
1399

- 1400 • Focuses on syntax, awkward phrasing, unnatural interrogative forms
1401

Score	Description
3	Fully natural and fluent; well-formed question
2	Mostly fluent; minor grammatical issues or slight awkwardness
1	Understandable but ungrammatical or clearly unnatural
0	Ungrammatical, confusing, or not a valid question

Score Description

3	The question fully matches the meaning of the original facts; no distortion or loss of critical details; context is complete and the answer remains uniquely correct.
2	The question is generally accurate but has minor factual ambiguity or slightly softer interpretation that could cause mild uncertainty while still pointing to the same answer.
1	Some key context or factual precision is missing or altered; relationships are weakened; the answer is still inferable but not strictly unique or reliable.
0	The question introduces clear factual errors, distorts the original meaning or relationships, or misleads so the answer could be wrong or invalid.

2) Factual Consistency: When given the answer, is the question factually consistent with the original sentence and triple (no distortion of meaning or relationships)? 1402
1403
1404
1405

- 1406 • **Semantic Accuracy:** Does the question preserve the intended meaning of the original sentence and triple? 1407
1408
- 1409 • **Relationship Integrity:** Are the logical relations (*e.g.*, cause/effect, association, identity) between subject, relation, and object kept intact? 1410
1411
- 1412 • **Context Preservation:** Does the question retain essential context needed to uniquely identify the correct answer (*e.g.*, disease type, anatomical site, conditions, time point, numerical thresholds)? Missing such context should lower the score. 1413
1414
1415
1416
1417

3) Linguistic Diversity: How well does the paraphrased question use different wording and structure from the original prompt? 1418
1419
1420

- 1421 • **What it checks:** Lexical and syntactic variation between the original and paraphrased versions 1422
- 1423 • **Focuses on:** Synonym usage, sentence structure changes, reformulation techniques 1424

4) Cross-Lang Factual Consistency: Can the evidence supporting the same fact be found or inferred 1425
1426

Score	Description
3	Excellent reformulation; uses different vocabulary and structure while maintaining meaning
2	Good variation; some different wording but follows similar structure or only changing the structure without introducing new vocabulary
1	Minimal variation; mostly replaces the blank with a question word
0	No meaningful reformulation; essentially the same as the original with a question mark

Score	Description
2	Direct match – same triple is clearly expressed in one sentence or consecutive sentences in the abstract in the target language.
1	Inferable without extra biomedical knowledge out of the given content – not in one sentence, but can be reasonably inferred from the paragraph as a whole.
0	Not supported or requires extra knowledge – the triple cannot be inferred from the abstract, or unrelated.

Score	Description
3	Highly plausible: Very convincing as an answer; can confuse even experts; fits subject, relation, domain well.
2	Moderately plausible: Makes sense in general; fits domain and context somewhat; can be ruled out by basic domain knowledge.
1	Barely plausible: Awkward or uncommon; easily ruled out by surface cues or common sense without any domain knowledge.
0	Implausible: Irrelevant, nonsensical, or grammatically incorrect; not a valid answer option.

Score	Description
3	Definitely wrong: contradicts or is not supported by the original sentence.
2	Likely wrong: but could be ambiguous or partially true given the original sentence.
1	Borderline: Possibly true or partially correct; ambiguous given the sentence.
0	Incorrectly labeled – This distractor is actually correct or the original answer given the original sentence.

1427 from the document in the target language?

1428 G.3 Distractor Quality Evaluation

1429 **Objective:** Evaluate the quality of three distractor
1430 options (incorrect candidates) accompanying the
1431 correct answer (object) in a multiple-choice setting
1432 derived from biomedical fill-in-the-blank prompts.
1433 Each distractor should be:

- 1434 • Plausible given the question
- 1435 • Incorrect (not the original object)
- 1436 • Relevant in context and domain

1437 **Input:**

- 1438 • Original Sentence
- 1439 • Triple
- 1440 • Fill-in-the-Blank Prompt
- 1441 • Answer options

1442 **Example:**

- 1443 • Original Sentence: EGFR is highly
1444 expressed in non-small cell lung
1445 carcinoma.
- 1446 • Triple: (*EGFR, is highly expressed in, non-small*
1447 *cell lung carcinoma*)
- 1448 • Prompt: EGFR is highly expressed in
1449 [BLANK].

1450 **Evaluation Criteria:** Need evaluations for both
1451 the cloze prompt and the paraphrased question.

1) Plausibility in Context: Is the distractor believable given the prompt and domain knowledge (biomedical)?

- **What it checks:** subject, relation, and expected answer type (should be the object)
- **Focuses on:** Be cautious of meaning shifts, incorrect substitutions, or role reversals.

2) Incorrectness: Is the distractor clearly incorrect given the original sentence and triple?

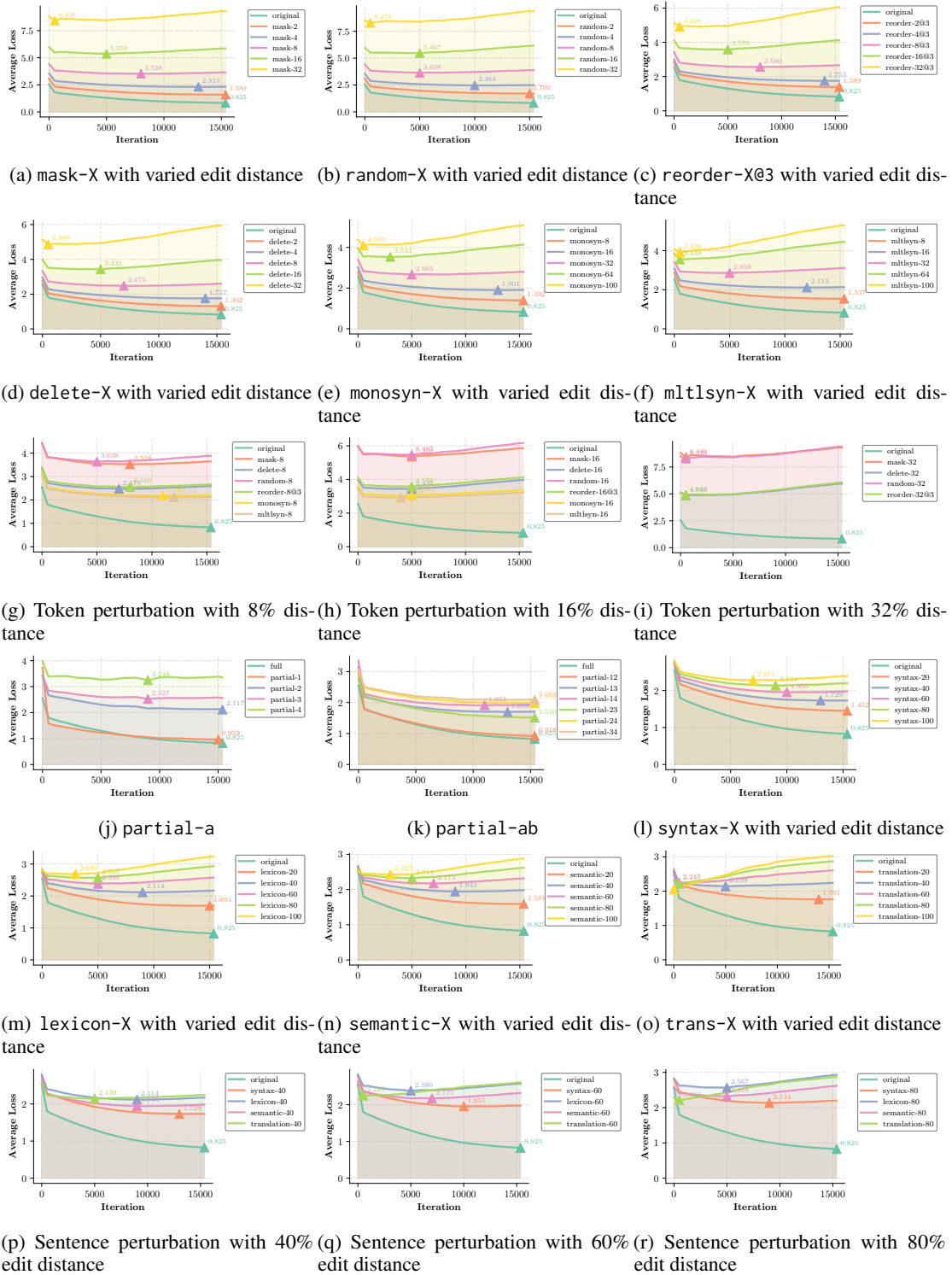


Figure 13: The loss dynamics over all perturbation patterns on Japanese sequences.

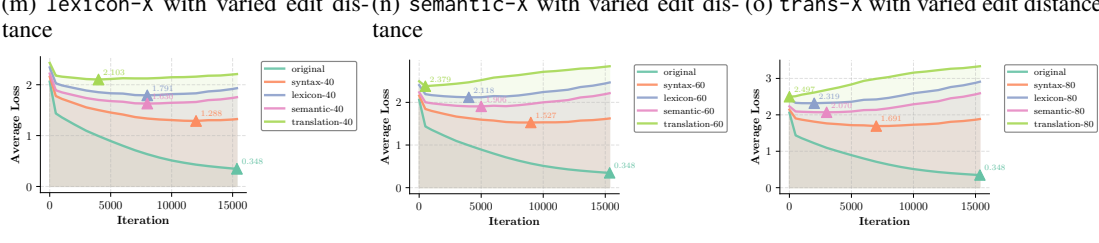
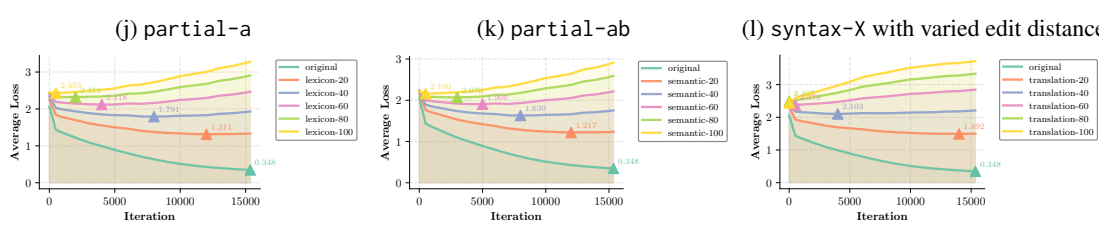
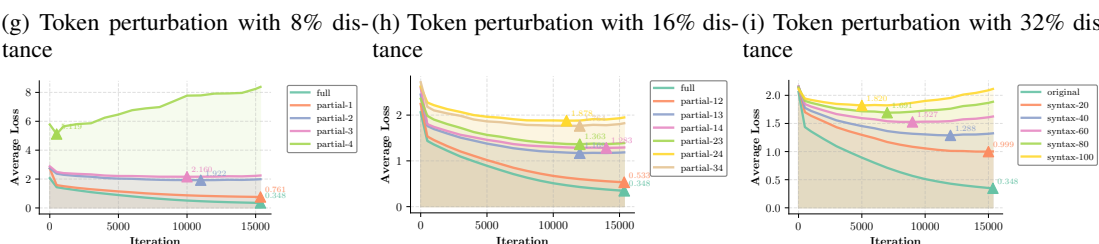
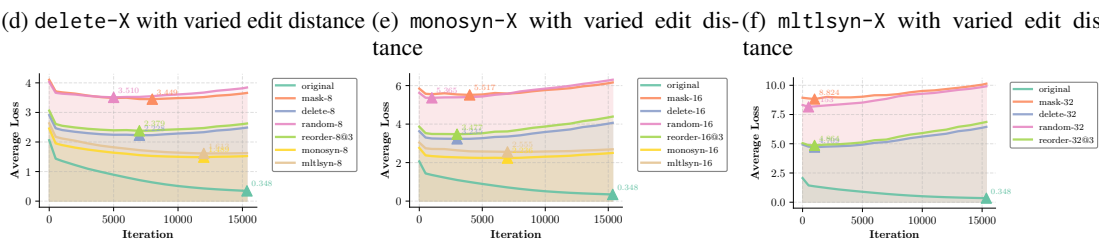
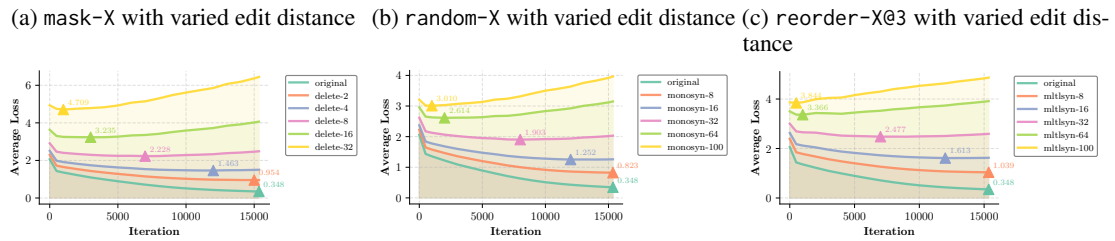
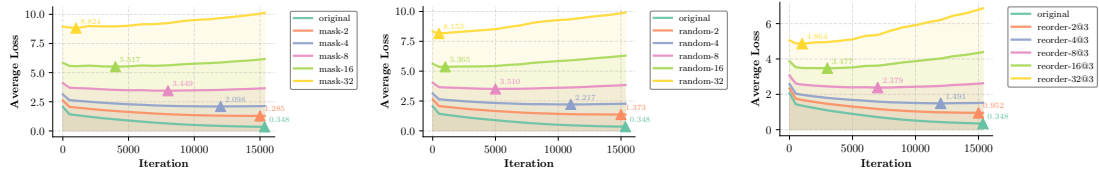


Figure 14: The loss dynamics over all perturbation patterns on English sequences.

Table 3: Samples of Japanese AdaXEval Dataset.

Sentence	cloze prompt	Paraphrase	Options	Answer ID
その結果, Ca-HApのカルシウムおよびリンの溶解は多くの有機酸水溶液中, pH4 6.5の範囲で下記の式に従うことがわかった(ただし, リンゴ酸, 酒石酸, クエン酸を除く)。	その結果, [BLANK]は多くの有機酸水溶液中, pH4~6.5の範囲で下記の式に従うことがわかった(ただし, リンゴ酸, 酒石酸, クエン酸を除く)。	多くの有機酸水溶液中, pH4~6.5の範囲で下記の式に従うことがわかったのは何の溶解ですか?	A. Ca-HApのマグネシウムおよび鉄の溶解 B. Ca-HApの亜鉛および銅の溶解 C. Ca-HApのカルシウムおよびリンの溶解 D. Ca-HApのナトリウムおよびカリウムの溶解	C
それとともに, β -lactamaseやaminoglycoside acetyltransferase (AAC)などによる抗菌薬の不活化やDNA gyraseの変化による抗菌薬親和性の減少, さらにbiofilm形成による抗菌薬の低浸透などが組み合わさり多剤耐性化する。	β -lactamaseやaminoglycoside acetyltransferase (AAC)による抗菌薬の不活化, DNA gyraseの変化による抗菌薬親和性の減少, さらにbiofilm形成による抗菌薬の低浸透などが組み合わさり[BLANK]する。	β -lactamaseやAACによる抗菌薬の不活化, DNA gyraseの変化, biofilm形成などが組み合わさることで生じる現象は何ですか?	A. 単剤耐性化 B. 多剤耐性化 C. 抗菌薬感受性の増加 D. 抗菌薬の副作用増強	B
悪性腫瘍に伴う血液凝固能亢進状態により脳卒中をきたす病態はTrousseau症候群として知られている。	悪性腫瘍に伴う血液凝固能亢進状態により脳卒中をきたす病態は[BLANK]として知られている。	悪性腫瘍に伴う血液凝固能亢進状態が原因で脳卒中を発症する病態は何と呼ばれていますか?	A. Goodpasture症候群 B. Trousseau症候群 C. 抗リン脂質抗体症候群 D. Lambert-Eaton症候群	B
根部よりの吸収は化合物の疎水性 (logP) と負の相関性を示した。	[BLANK]は化合物の疎水性 (logP) と負の相関性を示した。	化合物の疎水性 (logP) と負の相関性を示すのはどのような吸収ですか?	A. 根毛からの吸収 B. 茎部よりの吸収 C. 葉部よりの吸収 D. 根部よりの吸収	D
このことからp53はH2AXリン酸化には関与しておらず, 脱リン酸化やさらに下流の因子と関わっていると考えられる。	このことから[BLANK]はH2AXリン酸化には関与しておらず, 脱リン酸化やさらに下流の因子と関わっていると考えられる。	H2AXリン酸化に関与していないと考えられるタンパク質は何ですか?	A. DNA-PK B. CHK2 C. p53 D. ATM	C
Solitary fibrous tumor(SFT)は, 間葉系細胞由来の稀な腫瘍である。	[BLANK]は, 間葉系細胞由来の稀な腫瘍である。	間葉系細胞由来の稀な腫瘍として知られているのは何ですか?	A. 神経膠腫 (glioma) B. リンパ腫 (lymphoma) C. Solitary fibrous tumor (SFT) D. 扁平上皮癌 (squamous cell carcinoma)	C
1990年代半ばに開発された脱窒菌法により, 硝酸イオンの δ 15N, δ 18Oを微量で同時に測定できるようになった。	1990年代半ばに開発された[BLANK]により, 硝酸イオンの δ 15N, δ 18Oを微量で同時に測定できるようになった。	硝酸イオンの δ 15Nおよび δ 18Oを微量で同時に測定することを可能にした方法は何ですか?	A. 分光光度法 B. 脱窒菌法 C. イオンクロマトグラフィー法 D. ガスクロマトグラフィー質量分析法	B
遠隔転移は悪性腫瘍が全身化した状態で, 治療の原則は薬物療法である。	遠隔転移は悪性腫瘍が全身化した状態で, 治療の原則は[BLANK]である。	悪性腫瘍が全身化した状態である遠隔転移の治療の原則として用いられるのは何ですか?	A. 放射線療法 B. 免疫療法 C. 外科的切除 D. 薬物療法	D
門脈圧亢進症における消化管壁内粘膜下A-Vanastomosis (A-VA) 開大増加に伴う循環亢進状態に関しては, 食道胃静脈瘤・門脈圧亢進症性胃症の発症に直接関連する病態として, 多くの検討がなされている。	門脈圧亢進症における消化管壁内粘膜下A-Vanastomosisの開大増加は[BLANK]に直接関連する病態として, 多くの検討がなされている。	門脈圧亢進症における消化管壁内粘膜下A-Vanastomosisの開大増加が直接関連する病態として発症するのは何ですか?	A. 肝性脳症および肝腎症候群の発症 B. 胆道閉塞および胆石症の発症 C. 膈炎および十二指腸潰瘍の発症 D. 食道胃静脈瘤および門脈圧亢進症性胃症の発症	D
BLT1を介したシグナルはMyD88の遺伝子発現を誘導することで腸内細菌からの自然免疫シグナルを増強し, 形質細胞の細胞増殖を促進することで経口ワクチン抗原に対する抗原特異的IgA産生を促進する作用があることがわかった。	BLT1を介したシグナルはMyD88の遺伝子発現を誘導することで[BLANK]を増強し, 形質細胞の細胞増殖を促進することで経口ワクチン抗原に対する抗原特異的IgA産生を促進する作用があることがわかった。	BLT1を介したシグナルがMyD88の遺伝子発現を誘導することで増強するのは何のシグナルですか?	A. 腸管上皮細胞の増殖シグナル B. 抗原提示細胞による炎症性サイトカインシグナル C. 腸内細菌からの自然免疫シグナル D. ウイルス感染による獲得免疫シグナル	C

Table 4: Samples of English AdaXEval Dataset.

Sentence	cloze prompt	Paraphrase	Options	Answer ID
Fe2O3-SiO2 particles, which removes H2S and COS in hot coal gas, are prepared.	Fe2O3-SiO2 particles, which remove [BLANK] in hot coal gas, are prepared.	Which contaminants in hot coal gas are targeted for removal by Fe2O3-SiO2 particles?	A. NH3 and HCN B. SO2 and NOx C. H2S and COS D. CO2 and CH4	C
The antimicrobial activity of lomefloxacin against gram-positive bacteria was inferior to those of ofloxacin and gentamicin and comparable to that of chloramphenicol.	The antimicrobial activity of lomefloxacin against gram-positive bacteria was inferior to those of [BLANK] and comparable to that of chloramphenicol.	Against gram-positive bacteria, lomefloxacin shows lower antimicrobial activity than which other antibiotics?	A. ciprofloxacin and ampicillin B. ofloxacin and gentamicin C. vancomycin and clindamycin D. erythromycin and tetracycline	B
Basal cell adenoma is a rare type of salivary gland tumor.	[BLANK] is a rare type of salivary gland tumor.	Which rare type of tumor can occur in the salivary glands?	A. Adenoid cystic carcinoma B. Pleomorphic adenoma C. Basal cell adenoma D. Mucoepidermoid carcinoma	C
The common clinical signs were fever and hepatosplenomegaly.	The common clinical signs were [BLANK].	Which clinical signs are most frequently observed?	A. fever and hepatosplenomegaly B. jaundice and ascites C. rash and lymphadenopathy D. cough and chest pain	A
Phospholipids were found to be classified into three groups : (1) a lipid deactivating the glycolipid by strong hydrogen bond (phosphatidic acid analog), (2) a lipid likely to distribute the glycolipid rather homogeneously by weak hydrogen bond (phosphatidylglycerol and phosphatidylinositol analogs), and (3) a lipid enhancing the activity of a glycolipid by electrostatic effect (phosphatidylserine, phosphatidylcholine, and phosphatidylethanolamine analogs).	Phospholipids were found to be classified into three groups: (1) a lipid deactivating the glycolipid by strong hydrogen bond (phosphatidic acid analog), (2) a lipid likely to distribute the glycolipid rather homogeneously by weak hydrogen bond (phosphatidylglycerol and phosphatidylinositol analogs), and (3) [BLANK].	Which group of phospholipids enhances the activity of glycolipids by electrostatic effect?	A. lipids inhibiting glycolipid synthesis by covalent modification (ceramide analogs) B. lipids deactivating glycolipids by strong hydrogen bond (phosphatidic acid analog) C. lipids enhancing glycolipid activity by electrostatic effect (phosphatidylserine, phosphatidylcholine, and phosphatidylethanolamine analogs) D. lipids distributing glycolipids homogeneously by weak hydrogen bond (phosphatidylglycerol and phosphatidylinositol analogs)	C
Furthermore, the dried ARC, which was dehydrated in the presence of saccharides, can be recovered by dispersion of the powdered ARC in water.	The dried ARC, which was dehydrated in the presence of saccharides, can be recovered by [BLANK].	What process allows the recovery of dried ARC that was dehydrated with saccharides?	A. heating the powdered ARC in ethanol B. exposing the powdered ARC to ultraviolet light C. dispersion of the powdered ARC in water D. mixing the powdered ARC with organic solvents	C
About half of the reported cases of acanthosis nigricans are accompanied by various kinds of malignant neoplasms, mostly adenocarcinomas of the digestive system.	About half of the reported cases of acanthosis nigricans are accompanied by [BLANK].	What condition is present in about half of the reported cases of acanthosis nigricans?	A. various kinds of malignant neoplasms, mostly adenocarcinomas of the digestive system B. infectious diseases, such as tuberculosis and hepatitis C. autoimmune disorders, mainly lupus and rheumatoid arthritis D. benign skin tumors, primarily lipomas and fibromas	A
Suppressed cellular immunity, malignancy, diabetes mellitus and history of antibiotic usage are significant predisposing factors for the development of esophageal candidiasis.	[BLANK] are significant predisposing factors for the development of esophageal candidiasis.	Which conditions are considered important risk factors for developing esophageal candidiasis?	A. Chronic hypertension, obesity, hyperlipidemia, and smoking B. Suppressed cellular immunity, malignancy, diabetes mellitus, and history of antibiotic usage C. Asthma, seasonal allergies, eczema, and vitamin D deficiency D. Alcohol abuse, liver cirrhosis, renal failure, and hypothyroidism	B
Improved hydrogen sulfide removal is necessary for this apparatus to be applied to measurement of biogas produced by anaerobic digestion, since hydrogen sulfide influences catalysis.	Improved hydrogen sulfide removal is necessary for this apparatus to be applied to measurement of biogas produced by anaerobic digestion, since hydrogen sulfide influences [BLANK].	What process is affected by the presence of hydrogen sulfide in biogas measurement apparatus?	A. oxidation B. fermentation C. catalysis D. photosynthesis	C
The activity of ACC (1-aminocyclopropane-1-carboxylic acid) oxidase and the rate of ethylene production increased rapidly during fruit ripening at 20°C.	The activity of ACC (1-aminocyclopropane-1-carboxylic acid) oxidase and the rate of ethylene production increased rapidly during [BLANK].	During which process at 20°C do ACC oxidase activity and ethylene production rate increase rapidly?	A. fruit ripening at 20°C B. leaf senescence at 20°C C. flowering at 20°C D. seed germination at 20°C	A