Controlled Cloze-test Question Generation with Surrogate Models for IRT Assessment

Anonymous ACL submission

Abstract

Item difficulty plays a crucial role in 2 adaptive testing. However, few works have 3 focused on generating questions of varying 4 difficulty levels, especially for multiple-5 choice (MC) cloze tests. We propose 6 training pre-trained language models 7 (PLMs) as surrogate models to enable item 8 response theory (IRT) assessment, avoiding ٥ the need for human test subjects. We also 10 propose two strategies to control the 11 difficulty levels of both the gaps and the 12 distractors using ranking rules to reduce 13 invalid distractors. Experimentation on a 14 benchmark dataset demonstrates that our 15 proposed framework and methods can 16 effectively control and evaluate the 17 difficulty levels of MC cloze tests. 18

Introduction 19

27

1

Multiple-choice cloze tests are fill-in-the-blank 20 21 questions that assess reading comprehension and 22 overall language proficiency by requiring test 23 takers to select the correct missing words from 24 options. Table 1 gives an example test item 25 consisting of a stem with a gap to fill, a key or 26 answer, and three distractors.

	Stem:					
	I knelt and put my arms around the child. Then the tears					
	came, slowly at first, but soon she was her heart out					
	against my shoulder.					
	Options:					
	A. crying B. shouting C. drawing D. knocking					
_	Key: A Distractors: B C D					
28	Table 1: A question item of MC Cloze test.					
29						

30 ³¹ research because they are a common question ⁷¹ mimic Item Response Theory, bypassing the need 32 format on standardized language proficiency 72 for human test subjects. We will provide our dataset 33 exams such as TOEFL, TOEIC, IELTS, and 73 and codes upon request. 74

³⁴ college/high school entrance exams. In this paper, 35 we address the research questions of generating ³⁶ MC cloze test of different item difficulty levels.

Prior studies on cloze test question generation ³⁸ have concentrated largely on distractor generation, ³⁹ with the goal of reproducing distractors exactly 40 matching the benchmark datasets (Chung et al. 41 2020; Ren et al., 2021; Chiang et al. 2022; Wang et some 2023). Although studies have 42 al. 43 acknowledged the benefit of having distractors 44 with diverse difficulty levels (Yeung et al., 2019), 45 there has been minimal investigation into ⁴⁶ generating distractors with difficulty level different 47 from the benchmark.

Item difficulty plays a crucial role in adaptive 48 49 testing. It is a parameter that determines which 50 questions to present to a test taker and estimates ⁵¹ their proficiency level. Therefore, the difficulty of 52 each item should be known beforehand so 53 appropriate questions can be selected during the 54 test (Susanti et al. 2017). However, only a number 55 of works have focused on generating question ⁵⁶ items of various difficulty levels, for RC questions 57 (Gao et al. 2019a), C-test questions (Lee et al. 2019 ⁵⁸ and 2020), and MC cloze questions (Susanti et al. ⁵⁹ 2017). This research gap is largely due to the lack 60 of a reliable metric to evaluate the item difficulty ⁶¹ of generated questions. Most previous study relies 62 on human test takers and human annotation for 63 assessing the change of difficulty levels (Susantia 64 et al. 2017, Lee et al. 2019).

Our research has two main goals: (1) We propose 66 strategies to generate cloze-test questions by 67 controlling both the distractors and the gap, with 68 consideration for reducing invalid distractors. (2) ⁶⁹ We address the problem of objective and efficient MC cloze test questions have been a focus of 70 evaluation by using PLMs as subject surrogates to

			Factors to Control/Generate			Difficulty	
Related Research	Answer Type	Dataset	Distractor (Selection Method)	Gap (Generation Method)	Stem	Control (Evaluation Method)	Difficulty Level
Gao et al. 2019a	R. C.	SQuAD			\checkmark	Yes (RC system)	Item Level
Gao et al. 2019b	R. C.	RACE				None	
Chung et al. 2020	R. C.	RACE				None	
Qiu et al. 2020	R. C.	RACE	\checkmark			None	
Felice et al. 2022	Open Cloze	private		$\sqrt{(\text{Electra})}$		None	
Matsumori et al. 2023	Open Cloze	private		√ (gap score)		None	
Lee et al. 2019	C-test	Beiborn et al.2016		$\sqrt{(\text{prediction})}$		Yes (Human Subject)	Item Level
Lee et al. 2020	C-test	Beiborn et al.2016		√ (entropy)		Yes (MLP model)	Proficiency Level
Susantia et al. 2017	MC Cloze	TOEFL iBT	$\sqrt{(\text{feature-based})}$		\checkmark	Yes (Human subject)	Item Level
Yeung et al. 2019	MC Cloze	Chinese sentences	$\sqrt{(\text{BERT-based})}$ ranking)			None	
Ren and Zhu, 2021	MC Cloze	DGen	√ (featured-based L2R)			None	
Panda et al. 2022	MC Cloze	ESL lounge	$\sqrt{(\text{BERT-based and})}$			None	
Chiang et al. 2022	MC Cloze	CLOTH, DGen	$\sqrt{(\text{BERT-based and})}$			None	
Wang et al. 2023	MC Cloze	CLOTH, DGen	$\sqrt{(\text{Text2Text})}$			None	
Our Research	MC Cloze	CLOTH	√ (BERT-based and feature-based with validity rules)	√ (confidence- based entropy)		Yes (PLM- based IRT Assessment)	Item Level

75

76 77

Table 2: Recent Research on Question Generation for Language Proficiency Test

78 2 **Related Research**

The language proficiency test commonly adopts 79 80 cloze tests (open or multiple-choice), C-tests, and 81 reading comprehension (RC) to assess students' ⁸² language skills. Question Generation (OG) aims to 83 create natural and human-like questions from 84 diverse data sources. Research on MC cloze test 85 question generation primarily focuses on tasks 86 such as analyzing factors influencing item 87 difficulty (Susanti et al., 2017), distractor ⁸⁸ generation (Yeung et al., 2019; Ren and Zhu, 2021; 89 Chiang et al., 2022), and reducing invalid 90 distractors (Zesch and Melamud, 2014; Wojatzki et ⁹¹ al., 2016). Table 2 presents a comparative analysis 92 of recent studies on the automatic generation of 93 cloze test, RC, and C-test.

For MC cloze test, distractor generation 95 algorithms aim to identify plausible but incorrect 96 candidates for filling in blanks. Selection is based 97 on semantic proximity to the target word, measured 98 through methods like WordNet (Brown et al., 99 2005), thesauri (Smith et al., 2010), and word 100 embeddings similarity (Guo et al., 2016; Susanti et 101 al., 2015; Jiang and Lee, 2017). Recent studies 102 utilize confidence scores from BERT models

103 (Devlin et al. 2018) for ranking distractor 104 candidates, outperforming semantic similarity 105 methods in correlation with human judgment 106 (Yeung et al., 2019). Ren and Zhu (2021) apply 107 knowledge-based techniques to help generate 108 distractor candidates. Chiang et al. (2022) suggest 109 BERT-based methods as superior in distractor 110 generation. Their candidate selection relies on 111 confidence scores from pretrained language 112 models. Wang et al. (2023) propose a Text2Text 113 formulation using pseudo Kullback-Leibler 114 divergence, candidate augmentation and multi-task 115 training, enhancing performance in generating 116 distractors that align with benchmarks.

Item difficulty is crucial in adaptive testing, yet 118 few studies focus on generating items with diverse 119 difficulty levels different from standard benchmark 120 datasets. Furthermore, these works typically rely 121 on human test-taker evaluations (Susanti et al., 122 2017; Lee et al., 2019). A few studies used model 123 judgments in RC test (Gao et al., 2019) and C-test 124 (Lee et al., 2020). In related research on question 125 difficulty estimation, QA models are also proposed 126 to estimate difficulty through item response theory 127 (Benedetto, 2022).

Gap generation has been the focus in the context 158 3.1 128 129 of C-tests (Lee et al. 2019 and 2020). In open cloze 159 130 tests, Felice et al. (2022) recommend transformer 131 models and multi-objective learning for gap 132 prediction. Matsumori et al. (2023) propose a 133 masked language model approach with a gap score 134 metric for generating open cloze questions tailored 135 to specific target words. In contrast, research 136 addressing the control of difficulty levels by 137 modifying both distractors and gaps in multiple-138 choice cloze tests is lacking.

139 3 Methodology

Our research addresses key challenges in 171 140 141 generating MC cloze questions. We aim to produce 172 using BigBird (Zaheer et al. 2020) and Electra 142 questions with varying difficulty by managing both 173 (Clark et al. 2020) with different hyperparameters. 143 the gap and distractors, using ranking rules to eliminate invalid distractors. We also propose a 144 assessment framework to 145 PLM-based IRT 146 objectively evaluate item-level difficulty changes, alleviating reliance on human annotation. 147

148 149 training PLM-based models on benchmark data to 150 simulate test-takers; (2) designing strategies to 151 control difficulty by manipulating gaps and 152 distractors; (3) using PLM-based surrogate models 153 to take the modified tests and applying IRT to 154 evaluate difficulty changes.

155





Figure 1: Research Structure

with IRT **PLM-based** Assessment **Surrogate Models**

160 Calibrating test difficulty traditionally requires 161 trials with human subjects, which is time-162 consuming and costly. IRT is a framework to 163 estimate item difficulty unsupervised (Benedetto, 164 2022; Susanti et al., 2017). Previous work used 165 Reading Comprehension systems or MLP models 166 to evaluate change of predictions (Gao et al. 2019a, 167 Lee et al. 2020). We propose that predictions by 168 various PLMs with different settings can simulate 169 human test-taking for IRT without actual test-170 takers.

We fine-tune 12 PLM models on each dataset 174 Control strategies generate hard and easy versions 175 of each test fold. Trained surrogate models take 176 these versions, and their scores are aggregated 177 across folds. An IRT model fitted on the aggregated 178 scores for the original and modified tests evaluates As shown in Figure 1, our approach involves: (1) ¹⁷⁹ difficulty shifts between easy and hard versions by 180 modeling score distributions.

181 3.2 **Difficulty-controllable Question** 182 generation

For difficulty-controllable question generation, 183 184 We combine PLM-based confidence scores, 185 semantic similarity and edit distance metrics, and 186 validity rules to generate gaps and distractors at tunable difficulty levels. 187

Gap Difficulty Control: 189

Entropy has been studied as a proxy for gap 190 complexity in open cloze tests (Felice et al., 2022) and C-tests (Lee et al., 2020). We propose 192 leveraging pre-trained model confidence scores for 193 entropy estimation of candidate gaps, without 194 separate training (Figure 10, Appendix B). 195

Given a cloze stem, we identify candidate gap 196 words matching the POS tag of the original key. 197 We fine-tune a model like BERT to predict words 198 for each candidate gap. For each gap, we calculate 199 the Shannon entropy using the top K predictions 200

$$H(X) = -\sum_{i=1}^{n} p(x_i) \log_2 p(x_i)$$

where x_i is the i_{th} word predicted by BERT for 203 ²⁰⁴ the candidate gap, and $p(x_i)$ is BERT model's 205 confidence score for x_i .

201

202

We sort candidate gaps by entropy and select 256 Invalid Distractor Control 207 high entropy ones for hard questions and low 257 208 entropy for easy questions. Hard questions are 258 to correct answers, but selection based solely on 209 generated by selecting hard gaps and generating 259 semantic scores or PLM confidence may generate ²¹⁰ more difficult distractors for the selected gaps, and ²⁶⁰ invalid distractors that plausibly fit the gap. vice versa for easy questions. 211

213 Distractor Difficulty Factors:

214 215 Yeung et al. 2019, Ren and Zhu 2021, Chiang et al. 265 distractors by confidence scores produces 302 216 2022), we design three factors for distractor 266 items with 482 invalid distractors, 365 of which 217 generation: semantic similarity using word2vec 267 ranked higher than the correct answer. ²¹⁸ cosine similarity, syntactic similarity using 268 219 Levenshtein distance, and PLM confidence scores 269 research (Zesch and Melamud, 2014), we design 220 for gap prediction.

Confidence score

223 with our training data set, S be a cloze stem, V be 274 ratio are selected from 50 ranks after the answer. 225 a vocabulary list, A be the answer of S, and d_i be a 275 (3) For easy items, bottom two are chosen from word in V as a candidate distractor. We denote a ²⁷⁶ ranks 50-100 after the answer. 226 given stem S with the cloze blank filled in [Mask] 277 227 with $S_{\otimes [Mask]}$.

Confidence score C_i for d_i given by PLM is ²⁷⁹ 229 230 defined:

 $C_i = p(d_i | \mathbb{M}(S_{\otimes [Mask]})$

231 232

233

236

238

212

221

222

Semantic similarity

The semantic similarity S_i of the candidate 234 distractor and the answer is defined as: 235

 $S_i = CosineSimilarity(Embed(d_i), Embed(A))$ 237

where *Embed()* refers to Glove Embedding. 239

240 241

2

Levenshtein ratio

The Levenshtein ratio measures string similarity 242 on a scale from 0 to 1. It is defined as: 243

sum – ldist

44 Levenshtein_{ratio} =
$$\frac{Sur}{m}$$

sum ²⁴⁵ where *sum* is the total length of two strings, and 246 *ldist* is the weighted edit distance between two 247 strings based on Levenshtein distance (Levenshtein 248 et al. 1966). The Levenshtein distance counts 249 insertions, deletions, and substitutions to transform ²⁵⁰ one string into the other. The weighted distance is 251 calculated as:

Challenging distractors are semantically similar ²⁶¹ Previous work suggests context-sensitive lexical 262 inference rules can filter potentially appropriate 263 distractors (Zesch and Melamud, 2014). Our Inspired by previous work (Susantia et al. 2017, 264 analysis reveals that fine-tuning BERT and ranking

> Motivated by these observations and previous 270 validity rules:

(1) Valid distractors have lower confidence 271 272 ranking than answers. (2) For hard items, top two Formally, let M() be a PLM model finetuned 273 distractors by semantic similarity and Levenshtein

> Annotation and validity rule impact analysis are 278 in Appendix C.

280 Distractor Selection Strategy

With the three defined control factors and 281 282 validity rules, we design two strategies for 283 generating challenging or easy distractors. For 284 distractor generation, we use BERT to rank and 285 score all candidate words, then select the top 100 ²⁸⁶ ranks after the correct answer to form the distractor ²⁸⁷ candidate list, implementing the first validity rule.

The first strategy, Confidence-Ranking Control 288 289 (Figure 11, Appendix B), chooses the top 3 highest-290 confidence distractors from the candidate list for 291 difficult questions and the bottom 3 for easier 292 questions.

293 The second strategy, 3-Factor Ranking Control, 294 combines all three control factors – confidence 295 scores, word2vec embedding similarity, and ²⁹⁶ Levenshtein distance – along with validity rules 2 297 and 3 (Figure 12, Appendix B). This integrated ²⁹⁸ approach allows us to tune distractor difficulty.

299 4 **Experiment Design**

300 This section presents the experimentation details 301 for validating our proposed framework and 302 methods. Table 6 provides generation examples ³⁰³ referencing the original item shown in Table 1.





Table 3: Generated hard and easy items for the original item in Table

308 4.1 Dataset

309 ³¹⁰ generation (Table 1), with CLOTH¹ (Xie et al., ³⁴⁴ libraries. We apply the 1PL (also known as the 311 2017) and DGen (Ren & Zhu, 2021) being popular 345 Rash model) with default setting. This model 312 choices. DGen compiles science questions from 346 estimates a latent ability parameter for subjects and 313 diverse sources and levels, while CLOTH contains 347 a latent difficulty parameter for items, which fits 314 cloze-style English reading 315 questions for middle-school and high-school ³¹⁶ entrance exams. We selected CLOTH as it aligns ³⁴⁹ 5 ³¹⁷ closely with our goal of controlling item difficulty for adaptive testing. We strictly follow the "Terms 318 and Conditions" as listed on the download site. 319

We divided the CLOTH dataset into two sets 320 321 according to its two proficiency levels - CLOTH-322 M for middle school and CLOTH-H for high 323 school entrance exams. Each set was further 324 segmented into 5 folds. Within each fold, we split 325 the passages into stems. Stems comprised 326 consecutive sentences leading up to the first 327 [MASK] token (i.e. gap), ensuring sufficient 328 context surrounding the cloze deletion. Data 329 statistics is provided in Appendix A.

330 4.2 **Evaluation**

For each data fold, we trained 12 PLM models 331 332 using BigBird and Electra architectures, with 333 learning rates of 1e-4, 1e-5, and 3e-5, batch sizes 334 of 16 and 32, epoch of 1 and AdamW optimizer. We conducted experiments on a single NVIDIA 361 Table 4. Surrogate models' performance 335 336 Quadro RTX 8000 GPU. The control strategies 362 ³³⁷ were applied to the "Test" split. By concatenating ³⁶³ 338 the scores across all surrogate models and test 364 select 4 as middle school surrogates (Electra (1e-4, 339 folds, IRT models were then fitted to quantify 365 16), Electra (1e-4, 32), Bigbird (1e-4, 32), and 340 overall changes in test difficulty. We use the py-irt 366 Electra (3e-5, 32)) and 3 as high school surrogates

³⁴¹ library (Lalor and Rodriguez, 2023) as it leverages 342 PyTorch and GPU acceleration for faster and more Various datasets have been used for cloze test 343 scalable IRT modeling compared to existing comprehension 348 exactly what we intend to evaluate.

Results and Analysis

350 Surrogate Model Performance

Table 3 presents the surrogate models' average 352 accuracies across the five test data folds on the 353 original cloze items. Italic numbers indicate the 12 354 CLOTH-M surrogates' performances, while 355 underlined numbers show the 12 CLOTH-H 356 surrogates. The models exhibit a wide accuracy ³⁵⁷ range (0.42 to 0.81), demonstrating diverse 358 capabilities as artificial test takers for difficulty 359 modeling.

Proficiency	CLOT	TH-M	CLOT	ГН-Н
Model	BigBird	Electra	BigBird	Electra
1e-4, 16	0.4282	0.7106	0.4234	0.527
1e-4, 32	0.6691	0.7306	0.5671	<u>0.6601</u>
1e-5, 16	0.811	0.7613	0.7902	0.7119
1e-5, 32	0.8081	0.7602	<u>0.7974</u>	0.7102
3e-5, 16	0.6093	0.7558	<u>0.687</u>	<u>0.7008</u>
3e-5, 32	0.798	0.7615	0.7814	0.7072

To further analyze the surrogate models, we

¹ https://www.cs.cmu.edu/~glai1/data/cloth/

367 (Bigbird (1e-5, 16), Bigbird (1e-5, 32), Bigbird 415 CLOTH-M, retaining the original gap position 366 (3e-5, 32)). These models are trained similarly and 416 versus modifying it does not substantially impact 369 tested on both CLOTH-M and CLOTH-H. The 417 the efficacy of confidence ranking or 3-factor ³⁷⁰ table below compares the average accuracies, ⁴¹⁸ ranking (Fig. 4). The item difficulty distributions ³⁷¹ standard deviations, and utility ratios of the 12- ⁴¹⁹ remain relatively consistent. 372 surrogate sets and the middle and high school 420 373 surrogate subsets:

3	7	4	
-			_

	(CLOTH-N	Λ	(CLOTH-H	ł
Model	Avg.	Ctdy	Utility	Avg.	Ctdy	Utility
	Acc.	Stav	Ratio	Acc.	Stav	Ratio
12	0.717	0.104	73.9%	0.672	0.108	72.1%
4-mid	0.718	0.034	38.2%	0.615	0.072	52.2%
3-high	0.803	0.006	10.4%	0.79	0.007	10.9%

375 Table 5. Comparing surrogate models

376

The utility ratio is the percentage of test questions 377 remaining after excluding those answered correctly 378 or incorrectly by all test takers. The 4 middle 379 school surrogates perform better on CLOTH-M 380 and worse on CLOTH-H, while the 3 high school 381 surrogates substantially outscore them on CLOTH-382 383 M. The smaller standard deviations demonstrate ³⁸⁴ these sets represent distinct proficiency levels. The 385 12-surrogate sets achieve higher utility ratios (73.9%, 72.1%) than the middle and high school 386 ³⁸⁷ sets, and are retained for evaluating item difficulty 388 control given their better utility and diverse 389 performances to distinguish between stronger and weaker students. 390

391

Performance of Control Methods 392

393 difficult and easy items using the Confidence-394 ³⁹⁵ ranking algorithm and 3-Factor strategy. The red lines show IRT distributions for difficult generated 396 items, the blue lines for easy items, and the black 397 dotted lines mark the original test difficulty. 398

Both strategies systematically manipulated 399 400 cloze item difficulty. Across CLOTH-M and CLOTH-H, the strategies successfully generated 401 402 harder items (red distribution shift right) and easier 403 items (blue shift left) compared to the original test 404 items.

However, CLOTH-H exhibits a narrower spread 405 406 between high and low difficulty items, indicating greater efficacy adjusting difficulty for the lower 407 proficiency CLOTH-M rather than the advanced 408 CLOTH-H. 409

410

411 Effect of Gap Control

The effect of gap position control on difficulty 431 412 413 control differs for intermediate (CLOTH-M) 432 Figure 3. Change of IRT for CLOTH-H with Confidence-414 versus advanced (CLOTH-H) questions. For 433 Ranking (above) and 3-Factor Ranking Control (below)

In contrast, for CLOTH-H, retaining the original 421 gap positions leads to wider item difficulty 422 distributions when generating hard questions, as 423 shown by the rose dashed line and the red dashed 424 line in Figure 4. However, for generating easy 425 CLOTH-H items, controlling the gap position or 426 not produces similar difficulty distributions.













Figure 4: Comparing the Effect of Gap-control on Two Strategies for Two Datasets

437 Comparing Confidence-ranking and 3-Factor 438 Ranking on Generating Easy and Hard Items

Comparing the two control strategies, 3-factor ranking generates a slightly wider range of easy titem difficulties for both datasets as shown in Figure 5. Meanwhile, confidence-ranking method without gap control produces slightly wider datasets as states as shown in figure 6.



Figure 5: Confidence-ranking and 3-factor ranking w/
and w/o Gap control on generating easy items for
CLOTH-M (above) and CLOTH-H (below)





455 Best Combination Strategies and456 Advantage of Gap Control

We provide the box plot analysis on best combination control strategies for the two proficiency tests. Figures 7 and 8 show that the best strategy combination is the 3-Factor Ranking control without gap for easy question and confidence ranking control without gap for hard questions. Figure 9 shows Gap Control with 3-Factor ranking enhances easy CLOTH-M item generation over 3-Factor without gap control.

⁴⁶⁶ While maintaining the same mean difficulty, Gap ⁴⁸⁵ **6**

467 Control increases variability, indicating improved

469 fulfilling key needs when creating easy test 487 framework for assessing control of item-level 470 questions.

471



472 473 Figure 7: Best combination for CLOTH-M: 3-Factor 474 Ranking without Gap Control for Easy Items and 475 Confidence-Ranking without Gap Control for Hard Items 503 Comparatively, 3-Factor Ranking Control method 476 Generation.



478 Figure 8: Best combination for CLOTH-H: 3-Factor Ranking 479 without Gap Control for Easy Items and Confidence-Ranking 480 without Gap Control for Hard Items Generation.



482 Figure 9: Gap Control strategy increases item variability than 483 without gap control for 3-Factor Ranking method on easy 484 item generation for CLOTH-M.

Conclusions

468 ability to span multiple difficulty values - better 486 In this work, we proposed a novel evalutation 488 difficulty for MC Cloze test. By using diverse 489 pretrained models as surrogate test takers, we fitted 490 IRT distributions to quantify changes in difficulty avoiding reliance on human test subjects. 491

> We designed two strategies leveraging entropy, 492 493 semantic similarity, edit distance to manipulate both the gap position and distractor selection for difficultuy-controlled question generation. We 495 further implemented validity rules to reduce 496 generation of invalid distractors. 497

> Systematic experimentation shows: (1) The 499 advanced test (CLOTH-H) is more difficult to 500 control than intermediate test (CLOTH-M); (2) Gap control has a limited effect, yet increases item 501 variability for easy CLOTH-M generation; (3) 502 504 works better for easy items generation while 505 Confidence Ranking Control method exceeds at 506 hard item generation; (4) Validity rules reduce but 507 do not eliminate invalid distractors -- further study ⁵⁰⁸ into this challenge is desired.

Limitation 7 509

510 Our difficulty control methods worked better for ⁵¹¹ intermediate exam questions than advanced ones. 512 More research is needed to improve the methods' 513 ability to handle very complex test items. 514 Additionally, our techniques should be validated across other subject domains. Questions also 515 persist around optimizing validity methods to 517 avoid invalid distractors. We do not anticipate any 518 potential risks or ethical concerns arising from the 519 proposed framework for generating multiple-520 choice cloze test questions with controllable difficulty levels using pre-trained language models. 521 522

523 References

Benedetto, L. 2022. An assessment of recent 524 techniques for question difficulty estimation from 525 text. 526

Brown, J., Frishkoff, G., & Eskenazi, M. 2005. 527 Automatic question generation for vocabulary 528 assessment. In Proceedings of Human Language 529 Technology Conference and Conference on 530 Empirical Methods in 531 Natural Language Processing pages 819-826. 532

533 Chiang, S. H., Wang, S. C., & Fan, Y. C. 2022. CDGP: 586

- Automatic Cloze Distractor Generation based on 587 534
- 535 Pre-trained Language Model. In Findings of the
- 536 Association for Computational Linguistics: EMNLP 589
- 2022, pages 5835-5840. 537
- 538 Chung, H. L., Chan, Y. H., & Fan, Y. C. 2020. A BERT- 591
- based distractor generation scheme with multi-539
- training tasking and negative answer 540 strategies. arXiv preprint arXiv:2010.05384. 541
- 542 Clark, K., Luong, M. T., Le, Q. V., & Manning, C. D. 595
- 2020. Electra: Pretraining text encoders 543 as
- discriminators rather than generators. arXiv preprint 544 597 arXiv:2003.10555. 545
- 546 Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. 599
- 2018. Bert: Pre-training of deep bidirectional 600 547 transformers for language understanding. arXiv 548 preprint arXiv:1810.04805. 549
- 550 Felice, M., & Buttery, P. 2019. Entropy as a Proxy for 603
- Gap Complexity in Open Cloze Tests. 604 551
- In Proceedings of the International Conference on 605 552
- Recent Advances in Natural Language Processing 606 553 607
- (RANLP 2019), pages 323-327. 554

- Constructing open cloze tests using generation and 609 556 discrimination capabilities of transformers. arXiv 610 557
- preprint arXiv:2204.07237. 558
- Gao, Y., Bing, L., Chen, W., Lyu, M. R., & King, I. 612 559 2019(a). Difficulty controllable generation of 613 560 reading comprehension questions. In Proceedings 614 561 of the Twenty-Eighth International Joint Conference 615 562
- on Artificial Intelligence, pages 4968-4974. 563
- 564 Gao, Y., Bing, L., Chen, W., Lyu, M. R., & King, I. 617 2019(b). Generating Distractors for Reading 618 565 from Comprehension Questions Real 619 566 Examinations. In Proceedings of the AAAI 620 567 Conference on Artificial Intelligence, 33(01): 568 pages 6423-30. 569 622 570 Guo, Q., Kulkarni, C., Kittur, A., Bigham, J. P., & 623
- Brunskill, E. 2016, May. Questimator: Generating 624 571
- 572
- In IJCAI-16: Proceedings of the AAAI Twenty-Fifth 573 626
- International Joint Conference on Artificial 627 574
- Intelligence. 575
- 576 Jiang, S., & Lee, J. S. 2017. Distractor generation for 629
- 577
- the 12th Workshop on Innovative Use of NLP for 578 631
- Building Educational Applications pages 143-148. 579 632
- 580 Lalor, J. P., & Rodriguez, P. 2023. py-irt: A scalable 633
- item response theory library for python. INFORMS 634 581 Journal on Computing, 35(1), pages 5-13. 635 582
- 583 Lee, J. U., Schwan, E., & Meyer, C. M. 584
- In Proceedings of the 57th Annual Meeting of the 638 585

- Association for Computational Linguistics, pages 360-370.
- 588 Lee, J. U., Meyer, C. M., & Gurevych, I. 2020. Empowering active learning to jointly optimize system and user demands. arXiv preprint 590 arXiv:2005.04470.
- 592 Levenshtein, V. I. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In Soviet physics doklady Vol. 10, No. 8, pages 707-710

593

594

- 596 Matsumori, S., Okuoka, K., Shibata, R., Inoue, M., Fukuchi, Y., & Imai, M. 2023. Mask and Cloze: Automatic Open Cloze Question Generation Using 598 a Masked Language Model. IEEE Access, 11, pages 9835-9850.
- 601 Panda, S., Gomez, F. P., Flor, M., & Rozovskaya, A. 2022. Automatic Generation of Distractors for Fill-602 in-the-Blank Exercises with Round-Trip Neural Machine Translation. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop, pages 391-401.
- 555 Felice, M., Taslimipoor, S., & Buttery, P. 2022. 608 Qiu, Z., Wu, X., & Fan, W. 2020. Automatic distractor generation for multiple choice questions in standard tests. arXiv preprint arXiv:2011.13100.
 - 611 Ren, S., & Zhu, K. Q. 2021. Knowledge-driven distractor generation for cloze-style multiple choice questions. In Proceedings of the AAAI conference on artificial intelligence, Vol. 35, No. 5, pages 4339-4347
 - 616 Smith, S.; Avinesh, P.; and Kilgarriff, A. 2010. Gap-fill tests for language learners: Corpus-driven item generation. In Proceedings of ICON-2010: 8th International Conference on Natural Language Processing, pages 1-6.
 - 621 Susanti, Y. Iida, R. and Tokunaga, T. 2015. Automatic Generation of English Vocabulary Tests. In Proceedings of the 7th International Conference on Computer Supported Education, pages 77-87.
 - knowledge assessments for arbitrary topics. 625 Susanti, Y., Tokunaga, T., Nishikawa, H., & Obari, H. 2017. Controlling item difficulty for automatic vocabulary question generation. Research and practice in technology enhanced learning, 12(1):1-16.
 - chinese fill-in-the-blank items. In Proceedings of 630 Wang, H. J., Hsieh, K. Y., Yu, H. C., Tsou, J. C., Shih, Y. A., Huang, C. H., & Fan, Y. C. 2023. Distractor Generation based on Text2Text Language Models with Pseudo Kullback-Leibler Divergence Regulation. In Findings of the Association for Computational Linguistics: ACL 2023, pages 12477-12491.
 - 2019. Manipulating the Difficulty of C-Tests. 637 Wojatzki, M., Melamud, O., & Zesch, T. 2016. Bundled Gap Filling: A New Paradigm for Unambiguous Cloze Exercises. In Proceedings of

636

639

- the 11th Workshop on Innovative Use of NLP for 640 Building Educational Applications, pages 172–181. 641
- 642 Xie, Q., Lai, G., Dai, Z., & Hovy, E. 2018. Large-scale
- Cloze Test Dataset Created by Teachers. 643
- In Proceedings of the 2018 Conference on 644
- Empirical Methods in Natural Language 645

Processing, pages 2344-2356. 646

- 647 Yeung, C. Y., Lee, J. S., & Tsou, B. K. 2019. Difficulty-
- aware Distractor Generation for Gap-Fill Items. In 648
- Proceedings of the The 17th Annual Workshop of the 649
- Australasian Language Technology Association, 650
- pages 159-164. 651

652 Zaheer, M., Guruganesh, G., Dubey, K. A., Ainslie, J.,

Alberti, C., Ontanon, S., ... & Ahmed, A. 2020. Big 653

- bird: Transformers for longer sequences. Advances 654
- in Neural Information Processing Systems, 33, 655
- pages 17283-17297. 656

657 Zesch, T., & Melamud, O. 2014. Automatic Generation of Challenging Distractors Using Context-Sensitive 658

- Inference Rules. In Proceedings of the Ninth 672 659
- Workshop on Innovative Use of NLP for Building 660
- Educational Applications, pages 143-148. 661

662 A Data Statistics

⁶⁶³ Table 3 presents the number of items per split in our dataset. 664

Fold	Split	CLOTH-M	CLOTH-H
	Train	17123	42540
0	Validate	5678	14189
	Test	5669	14139
	Train	16975	42432
1	Validate	5757	14155
	Test	5738	14281
	Train	17011	42628
2	Validate	5680	14145
	Test	5779	14095
	Train	17094	42502
3	Validate	5733	14194
	Test	5643	14172
	Train	17077	42463
4	Validate	5752	14224
	Test	5641	14181

666 667

665

Table 6: Data Statistics

668 **B** Algorithms

669 Figures 8, 9, and 10 present the algorithms for

- 670 Gap Control, Confidence-Ranking Control, and 3
- 671 Factor Ranking Control respectively.

Algorithm 1: Gap generation **Input** : A sentence list with POS $A = [s_1, ..., s_n]$, Dict of POS

- numbers for answer words D_{POS} , Prediction model M, Numbers of model candidate K, Target level L
- **Output:** Target sentence list *T*
- $T \leftarrow$ Π;
- **2** for key in D_{POS} do $AllSentenceDict[key] \leftarrow [];$ 3
- 4 end
- 5 for sentence in A do
- for word in sentence \mathbf{do} 6
- $pos_{word} = pos(word);$ 7 if pos_{word} in D_{POS} then
- 8
- 10 // predict
- Sentence mask = CreateMaskSentence(sentence, word); $PredList = ModelPred(sentence_{mask}, K, M);$ // pre the top K score of $sentence_{mask}$ Using M 11
- shannon = CalShannon(PredList); // compute Shannon entropy of sentencemask
- 12 AllSentenceDict[posword].append([sentencemask, shannon]) $\dot{\mathbf{end}}$
- 13 end
- 14
- 15 end 16 for key in AllSentenceDict do
- 17
- $value = D_{POS}[key];$ if L is Hard then 18
- $T \leftarrow \text{Top } value \text{ candidate sentences sorted by } shannon \text{ in}$ 19 allSentenceDict[key]
- else if L is Easy then 20 $T \leftarrow \text{Last value candidate sentences sorted by shannon in}$ 21 allSentenceDict[key]

22 end 23 end

673 Figure 10: Gap Generation Algorithm

Algorithm 2: Distractor Generation with Confidence Ranking
Input : A sentence S with a cloze, Answer Word A, Prediction Model
M, Vocabulary List V , Numbers of Candidate K
Output: Difficult distractors D_{hard} , Easy distractors D_{easy}
$1 D \rightarrow D \leftarrow 0 D$

- 1 $D_{hard}, D_{easy} \leftarrow \{\}, \{\};$ 2 $ConfidenceList \leftarrow sorted(ModelPred(S, M, A, V));$
- // Sort
- 4 CandidateList = ConfidenceList[Position_A + 1 : Position_A + K];
- **5** $D_{hard} \leftarrow$ Top 3 candidates by *CandidateList*;
- 6 $D_{easy} \leftarrow \text{Last } 3 \text{ candidates by } CandidateList;$
- 7 return D_{hard} , D_{easy}

674

675 Figure 11: Distractor Generation with Confidence-Ranking Control 676

Algorithm 3: Distractor Generation with 3-Factor Ranking
Input : A sentence S with a cloze, Answer Word A, Prediction Model
M, Vocabulary List V , Numbers of Candidate K
Output: Difficult distractors D_{hard} , Easy distractors D_{easy}
1 $D_{hard}, D_{easy} \leftarrow \{\}, \{\};$
2 $ConfidenceList \leftarrow sorted(ModelPred(S, M, A, V));$ // Sort
vocabulary V based on the scores predicted by the model M
3 Position _A \leftarrow ConfidenceList.index(A);
4 CandidateList \leftarrow ConfidenceList[Position _A + 1 : Position _A + K];
5 $Similarity_{Glove}, Similarity_{Leven} \leftarrow [], [];$
6 for word in CandidateList do
7 $G \leftarrow \text{Calculate Glove similarity}(A, word);$
s $L \leftarrow \text{Calculate Leven similarity}(A, word);$
9 $Similarity_{Glove}$.append(G);
10 $Similarity_{Leven}$.append(L);
11 end
12 $GloveSimList_{hard}, GloveSimList_{easy} \leftarrow Similarity_{Glove}[Position_A + 1 :$
$Position_A + K/2$, $Similarity_{Glove}[Position_A + K/2 : Position_A + K]$;
13 LevenSimList _{hard} , LevenSimList _{easy} \leftarrow Similarity _{Leven} [Position _A + 1 :
$Position_A + K/2$, $Similarity_{Leven}[Position_A + K/2 : Position_A + K]$;

- 14 $D_{hard} \leftarrow$ Top 2 candidates by $GloveSimList_{hard}$, Top 1 candidate by LevenSimList_{hard}
- 15 $D_{easy} \leftarrow \text{Last 2 candidates by } GloveSimList_{easy}, \text{ Last 1 candidate by } LevenSimList_{easy};$

16 return D_{hard} , D_{easy}

Figure 12: Distractor Generation with 3-Factor Ranking Control 678

С Annotation for Invalid Distractor 679 Control 680

⁶⁸¹ We analyzed the issues of invalid distractors with 682 human evaluation. We recruited 9 college students 683 at the CET-6 English proficiency level as 684 annotators. The annotators work with our research

685 lab on regular basis and receive subsidy for their 727 (4) Distractors generated with Confidence 686 annotation work under supervision of our 728 Ranking Control: 687 administrative office. 729

The invalid distractors will most likely appear 730 Example #2: 689 when generating hard items. Using BERT's ⁶⁹⁰ confidence score ranking without validity control, ⁶⁹¹ we generated distractors for 4,575 items randomly 692 selected from the CLOTH-H dataset. Manual annotation identified 1,676 items as having at 694 least one invalid generated distractor (i.e., a 695 distractor that could fit as an answer in the gap). 736 (1) Original options: 696 As our control strategies involves ranking 737 697 distractors after the answer, we identified 302 ⁶⁹⁸ items to further test validity rules. Among the 906 699 distractors generated, 482 were annotated as ⁷⁰⁰ invalid, representing an invalidity ratio of 53.2%. ⁷⁴⁰ (3) Distractors generated with 3-Factor Ranking 701 After applying the Confidence-Ranking Control 741 Control: 702 method and 3-Factor Ranking Control method, 742 703 the ratios dropped to 20.3% and 17.3% 704 respectively (Table 7).

Strategy	Num. of Invalid Distractors	Ratio of Invalid Distractors
Confidence ranking w/o validity rules	482	53.2%
Confidence- ranking Control	184	20.3%
3-Factor Ranking Control	160	17.7%

705 Table 7. Manual annotation of 906 distractors 706 generated with confidence ranking w/o validity rules, 707 and our methods of Confidence-Ranking Control and 754 alternatives that contextually fit the gap as a 708 3-Factor Ranking Control

The following are examples of items with the 709 710 answer (bolded) and invalid 711 (italicized) generated by confidence ranking 712 without validity rules. The same item with 713 distractors generated using Confidence-ranking 759 When I began planning to move to Auckland to 714 Control and 3-Factor Ranking Control is also 760 study, my mother was worried about a lack of jobs 715 shown below:

716 Example #1:

717 I hope I did the right thing, Mom, Alice said. I saw 718 a cat, all bloody but alive. I [MASK] it to the vet's, 764 C. fears 719 and was asked to make payment immediately.

(1) Original options: 720

A. carried B. followed C. returned D. guided 721

722 (2) Distractors generated without control:

- A. carried B. took C. brought D. delivered 723
- 724 (3) Distractors generated with 3-Factor Ranking 725 Control:
- A. carried B. showed C. reported D. tried 726

A. carried B. transported C. hauled D. rode

731 [MASK] this surprised him very much, he went 732 through the paper twice, but was still not able to 733 find more than one mistake, so he sent for the 734 student to question him about his work after the 735 exam.

A. As B. For C. So D. Though

738 (2) Distractors generated without control: A. As B. Because C. Although D. Though 739

A. As B. Even C. Once D. Soon

743 (4) Distractors generated with Confidence 744 Ranking Control:

A. As B. Realizing C. Again D. Initially 745

747 D Instruction to Annotators for Invalid 748 Distractor Identification.

749 Instruction: You are given a set of multiple-750 choice cloze test questions, each with four options. 751 The correct answer is identified, along with three 752 generated distractor options. Please review the 753 choices and identify any "invalid distractors" -755 potentially correct response, rather than an 756 implausible one.

distractors 757 For example:

758 -----

⁷⁶¹ and cultural differences. Ignoring these I got 762 there in July 2010.

B. worries 763 A. concerns

D. considerations

765 -----

766 Here, the answer is "concerns". The generated 767 distractors include "worries". Both are "Concerns" 768 grammatically correct. fits the 769 semantic context only slightly better. Therefore, 770 in this case, "worries" is considered an "invalid 771 distractor".

772 Your annotation results will help assess the 773 efficacy of our difficulty-control strategies in 774 limiting invalid distractor generation for multiple 775 choice cloze tests.