

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

Anonymous ACL submission

Abstract

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

1 Introduction

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

Table 1: A question item of MC Cloze test.

MC cloze test questions have been a focus of research because they are a common question format on standardized language proficiency exams such as TOEFL, TOEIC, IELTS, and

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

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

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

Our research has two main goals: (1) We propose two strategies to generate cloze-test questions by controlling distractor difficulty, with consideration for reducing invalid distractors. (2) We address the problem of objective and efficient evaluation by using PLMs as subject surrogates to mimic Item Response Theory, bypassing the need for human test subjects. We will release our dataset and codes under CC-BY 4.0 license upon publication.

Related Research	Answer Type	Dataset	Factors to Control/Generate			Difficulty Control (Evaluation Method)	Difficulty Level
			Distractor (Selection Method)	Gap (Generation Method)	Stem		
Gao et al. 2019a	R. C.	SQuAD			√	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	√			None	
Felice et al. 2022	Open Cloze	private		√ (Electra)		None	
Matsumori et al. 2023	Open Cloze	private		√ (gap score)		None	
Lee et al. 2019	C-test	Beiborn et al.2016		√ (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	√ (feature-based)		√	Yes (Human subject)	Item Level
Yeung et al. 2019	MC Cloze	Chinese sentences	√ (BERT-based ranking)			None	
Ren and Zhu, 2021	MC Cloze	DGen	√ (feature-based L2R)			None	
Panda et al. 2022	MC Cloze	ESL lounge	√ (BERT-based and feature-based)			None	
Chiang et al. 2022	MC Cloze	CLOTH, DGen	√ (BERT-based and feature-based))			None	
Wang et al. 2023	MC Cloze	CLOTH, DGen	√ (Text2Text)			None	
Our Research	MC Cloze	CLOTH	√ (BERT-based and feature-based with validity rules)			Yes (PLM-based IRT Assessment)	Item Level

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

2 Related Research

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

For MC cloze test, distractor generation algorithms aim to identify plausible but incorrect candidates for filling in blanks. Selection is based on semantic proximity to the target word, measured through methods like WordNet (Brown et al., 2005), thesauri (Smith et al., 2010), and word embeddings similarity (Guo et al., 2016; Susanti et al., 2015; Jiang and Lee, 2017). Sakaguchi et al. (2013) introduce a discriminative approach for fill-

in-the-blank quiz generation for language learners. Recent studies utilize confidence scores from BERT models (Devlin et al. 2018) for ranking distractor candidates, outperforming semantic similarity methods in correlation with human judgment (Yeung et al., 2019). Ren and Zhu (2021) apply knowledge-based techniques to help generate distractor candidates. Chiang et al. (2022) suggest BERT-based methods as superior in distractor generation. Their candidate selection relies on confidence scores from pretrained language models. Wang et al. (2023) propose a Text2Text formulation using pseudo Kullback-Leibler divergence, candidate augmentation and multi-task training, enhancing performance in generating distractors that align with benchmarks.

Item difficulty is crucial in adaptive testing, yet few studies focus on generating items with diverse difficulty levels different from standard benchmark datasets. Furthermore, these works typically rely on human test-taker evaluations (Susanti et al., 2017; Lee et al., 2019). A few studies used model judgments in RC test (Gao et al., 2019) and C-test (Lee et al., 2020). Uto et al. (2023) propose difficulty-controllable neural question generation for reading comprehension using Item Response

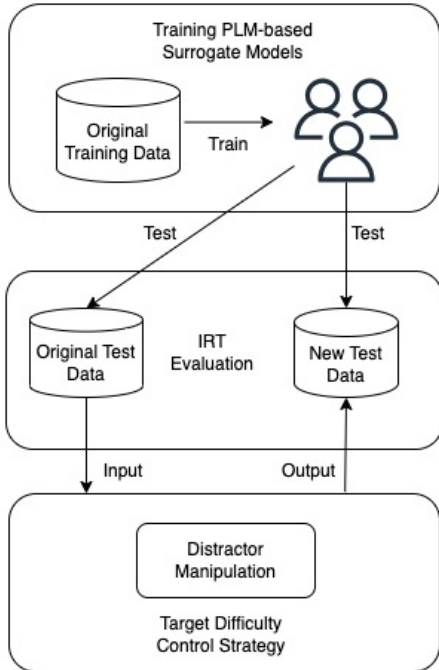
129 Theory. Ehara (2018) builds an English vocabulary
 130 knowledge dataset of Japanese English-as-a-
 131 Second-Language learners using crowdsourcing to
 132 assess learner proficiency. In related research on
 133 question difficulty estimation, QA models are also
 134 proposed to estimate difficulty through item
 135 response theory (Benedetto, 2022).

136 3 Methodology

137 Our research addresses key challenges in
 138 generating MC cloze questions. We aim to produce
 139 distractors with varying difficulty using ranking
 140 rules to eliminate invalid distractors. We also
 141 propose a PLM-based IRT assessment framework
 142 to objectively evaluate item-level difficulty
 143 changes, alleviating reliance on human annotation.

144 As shown in Figure 1, our approach involves: (1)
 145 training PLM-based models on benchmark data to
 146 simulate test-takers; (2) designing strategies to
 147 control difficulty by selecting distractors; (3) using
 148 PLM-based surrogate models to take the modified
 149 tests and applying IRT to evaluate difficulty
 150 changes.

151



152

153

Figure 1: Research Structure

154 3.1 IRT Assessment with PLM-based 155 Surrogate Models

156 Calibrating test difficulty traditionally requires
 157 trials with human subjects, which is time-
 158 consuming and costly. IRT is a framework to
 159 estimate item difficulty unsupervised (Benedetto,

160 2022; Susanti et al., 2017). Previous work used
 161 Reading Comprehension systems or MLP models
 162 to evaluate change of predictions (Gao et al. 2019a,
 163 Lee et al. 2020). We propose that predictions by
 164 various PLMs with different settings can simulate
 165 human test-taking for IRT without actual test-
 166 takers.

167 We fine-tune 12 PLM models on each dataset
 168 using BigBird (Zaheer et al. 2020) and Electra
 169 (Clark et al. 2020) with different hyperparameters.
 170 Control strategies generate hard and easy versions
 171 of each test fold. Trained surrogate models take
 172 these versions, and their scores are aggregated
 173 across folds. An IRT model fitted on the aggregated
 174 scores for the original and modified tests evaluates
 175 difficulty shifts between easy and hard versions by
 176 modeling score distributions.

177 3.2 Difficulty-controllable Question 178 generation

179 For difficulty-controllable question generation,
 180 we combine PLM-based confidence scores,
 181 semantic similarity and edit distance metrics, and
 182 validity rules to generate distractors at tunable
 183 difficulty levels.

184

185 Distractor Candidate List Generation

186 Let:

- 187 • $V = \{v_1, v_2, \dots, v_n\}$ be the proficiency
 188 vocabulary list, where v_i is the i -th word in the
 189 list;
- 190 • PTM be a pretrained model (e.g., BERT);
- 191 • $WP(v_i) = \{wp_1^i, wp_2^i, \dots, wp_{m_i}^i\}$ be the set of
 192 word pieces for word v_i after tokenization
 193 using PTM , where wp_j^i is the j -th word piece
 194 of v_i and m_i is the number of word pieces for
 195 v_i .
- 196 • Q be the item question.
- 197 • $R = \{r_1, r_2, \dots, r_k\}$ be the ranked list of the
 198 PTM word piece vocabulary, where r_i is the i -
 199 th ranked word piece and k is the total number
 200 of word pieces in the PTM vocabulary.

201 The algorithm to generate candidate distractors is
 202 as follows:

- 203 1. Tokenize each word $v_i \in V$ using PTM to
 204 obtain its word pieces $WP(v_i)$.
- 205 2. Predict the gap in question Q using PTM and
 206 obtain the ranked list R of the PTM word
 207 piece vocabulary.
- 208 3. For each word $v_i \in V$:

- 209 a. Calculate the rank of each word piece $wp_j^i \in$
 210 $WP(v_i)$ in the ranked list R . Let $rank(wp_j^i)$
 211 denote the rank of wp_j^i .
 212 b. Compute the mean rank of v_i 's word pieces:
 213 $rank(v_i) = \frac{1}{m_i} \sum_{j=1}^{m_i} rank(wp_j^i)$
 214 4. Sort the words in V based on their mean rank
 215 $rank(v_i)$ in ascending order to obtain the
 216 distractor candidate list $D = d_1, d_2, \dots, d_n$,
 217 where d_i is the i -th ranked word in the
 218 distractor candidate list.

219 Distractor Difficulty Factors:

220 Inspired by previous work (Susantia et al. 2017,
 221 Yeung et al. 2019, Ren and Zhu 2021, Chiang et al.
 222 2022), we design three factors for distractor
 223 generation: semantic similarity using word2vec
 224 cosine similarity, syntactic similarity using
 225 Levenshtein distance, and PLM confidence scores
 226 for gap prediction.

228 • Confidence score

229 Formally, let $\mathbb{M}()$ be a PLM model finetuned
 230 with our training data set, S be a cloze stem, V be a
 231 vocabulary list, A be the answer of S , and d_i be a
 232 word in V as a candidate distractor. We denote a
 233 given stem S with the cloze blank filled in $[Mask]$
 234 with $S_{\otimes[Mask]}$.

235 Confidence score C_i for d_i given by PLM is
 236 defined:

$$237 C_i = p(d_i | \mathbb{M}(S_{\otimes[Mask]}))$$

238 • Semantic similarity

239 The semantic similarity S_i of the candidate
 240 distractor and the answer is defined as:

$$241 S_i = \text{CosineSimilarity}(\text{Embed}(d_i), \text{Embed}(A))$$

242 where $\text{Embed}()$ refers to Glove Embedding.

243 • Levenshtein ratio

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

$$246 \text{Levenshtein}_{ratio} = \frac{\text{sum} - \text{ldist}}{\text{sum}}$$

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

$$254 \text{ldist} = \text{Num}(\text{INSERT}) + \text{Num}(\text{DELETE}) + 2$$

$$255 * \text{Num}(\text{REPLACE})$$

256 Distractor Selection Strategy with Validity 257 Control

258 In the context of distractor selection for
 259 generating challenging or easy test items, we
 260 introduce two strategies: Confidence-Ranking
 261 Control and 3-Factor Ranking Control. These
 262 strategies aim to select distractors that are
 263 semantically similar to the correct answer while
 264 ensuring their validity.

265 Let $V = \{v_1, v_2, \dots, v_n\}$ denote the vocabulary
 266 list, where v_i represents the i -th word in the list.
 267 The correct answer for a given test item is denoted
 268 by a , and the set of distractors is represented by
 269 $D = \{d_1, d_2, \dots, d_m\}$, where d_i is the i -th
 270 distractor.

271 We define $rank_{conf}(v_i)$, $rank_{sem}(v_i)$, and
 272 $rank_{lev}(v_i)$ as the ranks of word v_i in the
 273 vocabulary list V based on BERT confidence
 274 scores, semantic similarity scores, and Levenshtein
 275 ratio scores, respectively. Additionally, let $B =$
 276 $\{b_1, b_2, \dots, b_l\}$ be the set of benchmark distractors,
 277 where b_i is the i -th benchmark distractor.

278 Previous work proposes context-sensitive
 279 lexical inference rules to filter distractors (Zesch
 280 and Melamud, 2014). Our analysis as in Appendix
 281 C reveals fine-tuned BERT often ranks invalid
 282 distractors higher than correct answers. Motivated
 283 by this and prior research (Zesch and Melamud,
 284 2014), we design rules to reduce invalid distractors:

$$285 \forall d_i \in D, rank_{conf}(d_i) > rank_{conf}(a)$$

286 This rule guarantees that the selected distractors
 287 have lower confidence ranking positions than the
 288 correct answer.

289 The Confidence-Ranking Control Strategy
 290 selects distractors based on their BERT confidence
 291 scores. For hard distractors, we choose the top three
 292 distractors after the answer, as defined by:

$$293 D_h = \{v_i \in V \mid rank_{conf}(a) < rank_{conf}(v_i) \leq$$

$$294 rank_{conf}(a) + 3\}$$

295 For easy distractors, we select the first three
 296 distractors after the last ranked benchmark
 297 distractor, as defined by:

$$298 D_e = \{v_i \in V \mid rank_{conf}(b_l) < rank_{conf}(v_i)$$

$$299 \leq rank_{conf}(b_l) + 3\}$$

300 The 3-Factor Ranking Control Strategy
 301 considers semantic similarity, Levenshtein ratio,
 302 and BERT confidence scores to select distractors.

Let K be a predefined constant that determines the range of ranks to consider after the answer or the lowest confidence-ranked benchmark distractor.

For hard distractors, we select the top two distractors by semantic similarity and the top distractor by Levenshtein ratio from the K ranks after the answer, as defined by:

$$D_{h1} = \{v_i \in V \mid \text{rank}_{\text{conf}}(a) < \text{rank}_{\text{conf}}(v_i) \leq \text{rank}_{\text{conf}}(a) + K, \text{rank}_{\text{sem}}(v_i) \leq 2\}$$

$$D_{h2} = \{v_i \in V \mid \text{rank}_{\text{conf}}(a) < \text{rank}_{\text{conf}}(v_i) \leq \text{rank}_{\text{conf}}(a) + K, \text{rank}_{\text{lev}}(v_i) = 1\}$$

$$D_h = D_{h1} \cup D_{h2}$$

For easy distractors, we select the top two distractors by semantic similarity and the top distractor by Levenshtein ratio from the K ranks after the lowest confidence-ranked benchmark distractor, as defined by:

$$D_{e1} = \{v_i \in V \mid \text{rank}_{\text{conf}}(b_l) < \text{rank}_{\text{conf}}(v_i) \leq \text{rank}_{\text{conf}}(b_l) + K, \text{rank}_{\text{sem}}(v_i) \leq 2\}$$

$$D_{e2} = \{v_i \in V \mid \text{rank}_{\text{conf}}(b_l) < \text{rank}_{\text{conf}}(v_i) \leq \text{rank}_{\text{conf}}(b_l) + K, \text{rank}_{\text{lev}}(v_i) = 1\}$$

$$D_e = D_{e1} \cup D_{e2}$$

4 Experiment Design

This section presents the experimentation details. Table 3 provides generation examples referencing the original item shown in Table 1.

4.1 Dataset

Various datasets have been used for cloze test generation (Table 1), with CLOTH¹ (Xie et al., 2017) and DGen (Ren & Zhu, 2021) being popular choices. DGen compiles science questions from diverse sources and levels, while CLOTH contains cloze-style English reading comprehension questions for middle-school and high-school entrance exams. We selected CLOTH as it aligns closely with our goal of controlling item difficulty

for adaptive testing. We strictly follow the ‘‘Terms and Conditions’’ as listed on the download site.

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

4.2 Evaluation

We conducted experiments on a single NVIDIA Quadro RTX 8000 GPU. The control strategies were applied to the ‘‘Test’’ split. By concatenating the scores across all surrogate models and test folds, IRT models were then fitted to quantify overall changes in test difficulty. We use the py-irt library (Lalor and Rodriguez, 2023) as it leverages PyTorch and GPU acceleration for faster and more scalable IRT modeling compared to existing libraries. We apply the 1PL (also known as the Rash model) with default setting. This model estimates a latent ability parameter for subjects and a latent difficulty parameter for items, which fits exactly what we intend to evaluate.

5 Results and Analysis

Surrogate Model Performance

Table 4 presents the surrogate models’ average accuracies across the five test data folds on the original cloze items. Italic numbers indicate the 12 CLOTH-M surrogates’ performances, while underlined numbers show the 12 CLOTH-H surrogates. The models exhibit a wide accuracy range (0.42 to 0.81), demonstrating diverse capabilities as artificial test takers for difficulty modeling.

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

Table 4. Surrogate models’ performance

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

	Distractor generation w/ Confidence-Ranking Control	Distractor generation w/ 3-Factor Ranking Control
Hard	(I) <hr/> 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. <hr/> Options: A. crying B. sobbing C. pouring D. weeping <hr/> Key: A Distractors: B C D	(II) <hr/> 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. <hr/> Options: A. crying B. screaming C. cried D. crushed <hr/> Key: A Distractors: B C D
Easy	(III) <hr/> 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. <hr/> Options: A. crying B. counting C. shouting D. booming <hr/> Key: A Distractors: B C D	(IV) <hr/> 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. <hr/> Options: A. crying B. owing C. caves D. sobbed <hr/> Key: A Distractors: B C D

399

400

401

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

402

403

404

405

406

407

408

409

410

411

412

413

To further analyze the surrogate models, we select 4 as middle school surrogates (Electra (1e-4, 16), Electra (1e-4, 32), Bigbird (1e-4, 32), and Electra (3e-5, 32)) and 3 as high school surrogates (Bigbird (1e-5, 16), Bigbird (1e-5, 32), Bigbird (3e-5, 32)). These models are trained similarly and tested on both CLOTH-M and CLOTH-H. The table below compares the average accuracies, standard deviations, and utility ratios of the 12-surrogate sets and the middle and high school surrogate subsets:

Model	CLOTH-M			CLOTH-H		
	Avg. Acc.	Stdv	Utility Ratio	Avg. Acc.	Stdv	Utility 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%

414

415

416

Table 5. Comparing surrogate models

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

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

Performance of Control Methods

433

434

435

436

437

438

439

440

441

442

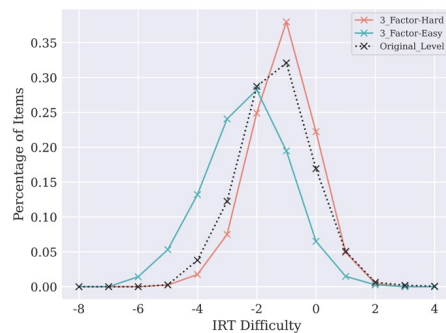
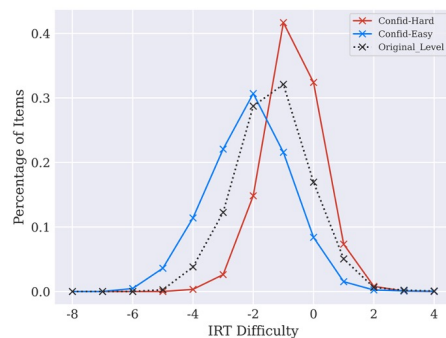
443

444

445

446

Figures 2 and 3 present the effect of generating difficult and easy items using the Confidence-ranking algorithm and 3-Factor strategy. The red lines show IRT distributions for difficult generated items, the blue lines for easy items, and the black dotted lines mark the original test difficulty. Both strategies systematically manipulated cloze item difficulty. Across CLOTH-M and CLOTH-H, the strategies successfully generated harder items (red distribution shift right) and easier items (blue shift left) compared to the original test items.

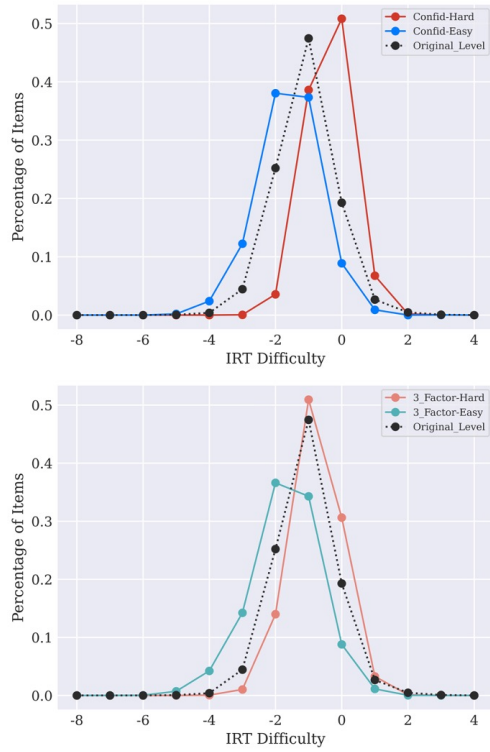


447

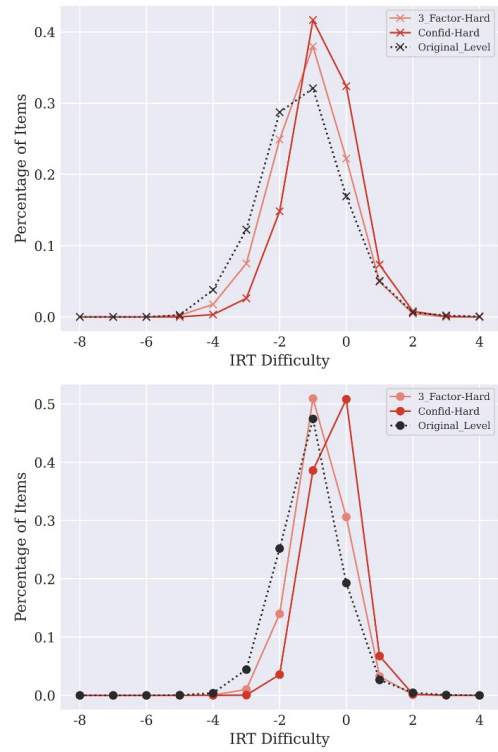
448

449

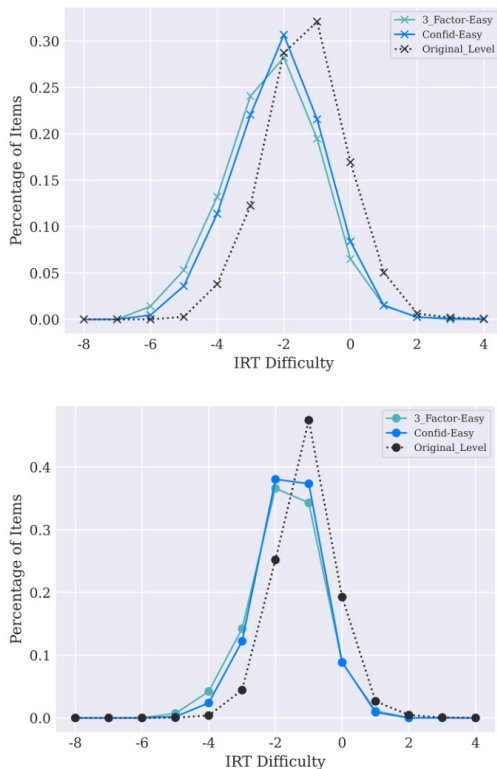
Figure 2. Change of IRT for CLOTH-M with Confidence-Ranking (above) and 3-Factor Ranking Control (below)



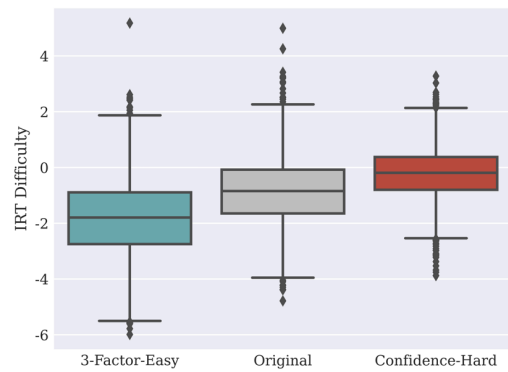
450
451 Figure 3. Change of IRT for CLOTH-H with Confidence-
452 Ranking (above) and 3-Factor Ranking Control (below)
453



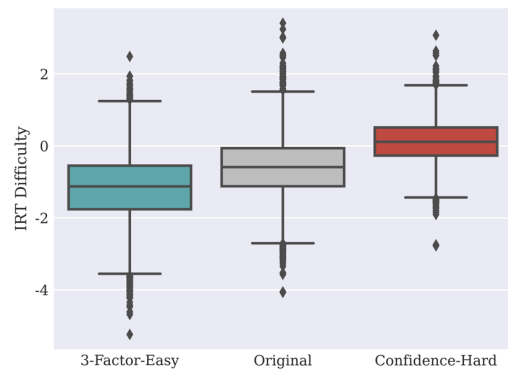
458
459 Figure 5: Confidence-Ranking and 3-Factor Ranking on
460 generating hard items for CLOTH-M (above) and CLOTH-H
461 (below)
462



454
455 Figure 4: Confidence-ranking and 3-factor ranking on
456 generating easy items for CLOTH-M (above) and
457 CLOTH-H (below)



463
464 Figure 6: Best combination for CLOTH-M: 3-Factor
465 Ranking for Easy Items and Confidence-Ranking for Hard
466 Items Generation.



467
468 Figure 7: Best combination for CLOTH-H: 3-Factor Ranking
469 for Easy Items and Confidence-Ranking for Hard Items
470 Generation.

471 Comparing Confidence-ranking and 3-Factor 472 Ranking on Generating Easy and Hard Items

473 Comparing the two control strategies, 3-factor
474 ranking generates a slightly wider range of easy
475 item difficulties for both datasets as shown in
476 Figure 4. Meanwhile, confidence-ranking method
477 produces slightly wider distributions of hard items
478 for both datasets as evident in Figure 5.

480 Best Combination Strategies

481 We provide the box plot analysis on best
482 combination control strategies for the two
483 proficiency tests. Figures 6 and 7 show that the best
484 strategy combination is the 3-Factor Ranking
485 control for easy question and confidence ranking
486 control for hard questions

487 6 Human Evaluation

488 We hired three college instructors from the
489 Faculty of English Language and Culture to
490 perform a human evaluation of our generated
491 questions. Two of them hold Ph.D. degrees in
492 Applied Linguistics and one with two Master
493 Degrees in Applied Linguistics and Computer-
494 assisted Language Teaching respectively. They
495 were compensated with a standard hourly rate for
496 college instructors, covering four hours of training
497 and four hours of annotation. Their payment was
498 supervised and approved by the College
499 Administration Office.

500 We randomly sampled 10 CLOTH-M articles
501 containing 100 cloze items and 6 CLOTH-H
502 articles containing 102 cloze items, comprising a
503 total of 100 easy items and 102 hard items. The
504 number of easy and hard items was equally split
505 within both the CLOTH-M and CLOTH-H
506 datasets. For each cloze item in each article, we
507 randomly paired the benchmark options with the
508 easy or hard items for the annotators to score their
509 relative difficulty. We then curated their responses
510 and transferred the scores as follows: the
511 benchmark was assigned a score of 2, the higher
512 annotated score was assigned a 3, and the lower
513 annotation was assigned a 1. The average results
514 are presented below.

515 The human evaluation results demonstrate that
516 the generated hard cloze items are consistently
517 perceived as more difficult than the easy items
518 across both the CLOTH-M and CLOTH-H datasets.
519 The total average scores for easy and hard items are
520 1.41 and 2.13, respectively, indicating a clear
521 distinction in difficulty levels. This trend is also

522 observed within each dataset, with the CLOTH-H
523 dataset showing a more pronounced difference
524 between easy and hard item scores compared to the
525 CLOTH-M dataset. Despite some variations
526 among the annotators, particularly in the CLOTH-
527 H dataset, the overall results confirm the
528 effectiveness of the generation process in creating
529 cloze items with varying levels of difficulty.

Total Avg.	Ann. #1	Ann. #2	Ann. #3	Avg.
Easy	1.22	1.12	1.9	1.41
Hard	1.93	2.04	2.42	2.13

530 (a) Average annotation scores for the complete set

CLOTH-M	Ann. #1	Ann. #2	Ann. #3	Avg.
Easy	1.28	1.16	1.36	1.27
Hard	1.94	2.28	1.92	2.05

531 (b) Average annotation scores for CLOTH-M

CLOTH-H	Ann. #1	Ann. #2	Ann. #3	Avg.
Easy	1.16	1.08	2.44	1.56
Hard	1.92	1.81	2.92	2.22

532 (c) Average annotation scores for CLOTH-H

533 Table 6: Human Annotation Results

537 7 Conclusions

538 In this work, we propose a novel evaluation
539 framework for assessing the control of item-level
540 difficulty in multiple-choice cloze tests. By
541 utilizing diverse pre-trained language models as
542 surrogate test-takers, we fit Item Response Theory
543 (IRT) distributions to quantify changes in difficulty,
544 avoiding reliance on human subjects.

545 We design two strategies leveraging confidence
546 scores, semantic similarity, and edit distance to
547 control distractor selection for generating questions
548 with controlled difficulty levels. To reduce the
549 generation of invalid distractors, we implement
550 validity rules based on prior research.

551 Systematic experimentation shows that (1) the
552 advanced test (CLOTH-H) is more challenging to
553 control than the intermediate test (CLOTH-M); (2)
554 the 3-Factor Ranking Control method is more
555 effective for generating easy items, while the
556 Confidence Ranking Control method excels at
557 generating hard items; (3) validity rules help
558 reduce invalid distractors but do not eliminate them
559 entirely, indicating a need for further research.

560 Our framework provides a promising approach
561 for generating multiple-choice cloze questions with
562 controllable difficulty levels, enabling more
563 effective adaptive testing.

564 8 Limitation

565 Our study has several limitations that provide
566 opportunities for future research. First, we
567 acknowledge that our sample size of 12 surrogate
568 models is smaller than conventionally
569 recommended for Item Response Theory (IRT)
570 analysis (Sireci, 1992). While our work explores
571 the potential of using IRT with smaller samples,
572 further investigation is needed to determine the
573 optimal number of surrogate models and how well
574 they mimic human test-takers.

575 Second, our study does not directly address the
576 question of how well the surrogate models simulate
577 actual test-takers, such as middle school and high
578 school students taking entrance exams.
579 Investigating this would require access to test
580 scores from human subjects, which raises
581 important considerations about data collection,
582 curation, and test-taker anonymity that are beyond
583 the scope of our current work.

584 Third, we recognize that IRT evaluates question
585 difficulty based solely on test-taker responses
586 without considering question content. Future
587 research should explore how different parameter
588 settings impact question quality assessment.

589 Finally, we opted for the 1PL (Rasch) model due
590 to its simplicity and suitability for smaller sample
591 sizes, as well as its alignment with our main
592 objective of evaluating the effectiveness of our
593 methods in controlling question difficulty.
594 However, we acknowledge that the 2PL model may
595 provide a more comprehensive understanding of
596 item properties, albeit at the cost of requiring larger
597 sample sizes for accurate parameter estimation.

598 Despite these limitations, we believe our work
599 provides a valuable starting point for future
600 research on using PLMs as surrogate test-takers
601 and applying IRT to assess the difficulty of
602 automatically generated questions. We hope our
603 study will inspire further investigations into these
604 important areas.

605 References

606 Benedetto, L. 2022. An assessment of recent
607 techniques for question difficulty estimation from
608 text.

609 Brown, J., Frishkoff, G., & Eskenazi, M. 2005.
610 Automatic question generation for vocabulary
611 assessment. In *Proceedings of Human Language
612 Technology Conference and Conference on
613 Empirical Methods in Natural Language
614 Processing* pages 819-826.

615 Chiang, S. H., Wang, S. C., & Fan, Y. C. 2022. CDGP:
616 Automatic Cloze Distractor Generation based on
617 Pre-trained Language Model. In *Findings of the
618 Association for Computational Linguistics: EMNLP
619 2022*, pages 5835–5840.

620 Chung, H. L., Chan, Y. H., & Fan, Y. C. 2020. A BERT-
621 based distractor generation scheme with multi-
622 tasking and negative answer training
623 strategies. *arXiv preprint arXiv:2010.05384*.

624 Clark, K., Luong, M. T., Le, Q. V., & Manning, C. D.
625 2020. Electra: Pretraining text encoders as
626 discriminators rather than generators. *arXiv preprint
627 arXiv:2003.10555*.

628 Devlin, J., Chang, M. W., Lee, K., & Toutanova, K.
629 2018. Bert: Pre-training of deep bidirectional
630 transformers for language understanding. *arXiv
631 preprint arXiv:1810.04805*.

632 Ehara, Y. (2018, May). Building an English vocabulary
633 knowledge dataset of Japanese English-as-a-
634 second-language learners using crowdsourcing.
635 In *Proceedings of the Eleventh International
636 Conference on Language Resources and Evaluation
637 (LREC 2018)*.

638 Felice, M., & Buttery, P. 2019. Entropy as a Proxy for
639 Gap Complexity in Open Cloze Tests.
640 In *Proceedings of the International Conference on
641 Recent Advances in Natural Language Processing
642 (RANLP 2019)*, pages 323–327.

643 Felice, M., Taslimipoor, S., & Buttery, P. 2022.
644 Constructing open cloze tests using generation and
645 discrimination capabilities of transformers. *arXiv
646 preprint arXiv:2204.07237*.

647 Gao, Y., Bing, L., Chen, W., Lyu, M. R., & King, I.
648 2019(a). Difficulty controllable generation of
649 reading comprehension questions. In *Proceedings
650 of the Twenty-Eighth International Joint Conference
651 on Artificial Intelligence*, pages 4968–4974.

652 Gao, Y., Bing, L., Chen, W., Lyu, M. R., & King, I.
653 2019(b). Generating Distractors for Reading
654 Comprehension Questions from Real
655 Examinations. In *Proceedings of the AAAI
656 Conference on Artificial Intelligence*, 33(01):
657 pages 6423–30.

658 Guo, Q., Kulkarni, C., Kittur, A., Bigham, J. P., &
659 Brunskill, E. 2016, May. Questimator: Generating
660 knowledge assessments for arbitrary topics.
661 In *IJCAI-16: Proceedings of the AAAI Twenty-Fifth
662 International Joint Conference on Artificial
663 Intelligence*.

664 Jiang, S., & Lee, J. S. 2017. Distractor generation for
665 chinese fill-in-the-blank items. In *Proceedings of
666 the 12th Workshop on Innovative Use of NLP for
667 Building Educational Applications* pages 143-148.

- 668 Lalor, J. P., & Rodriguez, P. 2023. py-irt: A scalable
669 item response theory library for python. *INFORMS*
670 *Journal on Computing*, 35(1), pages 5-13.
- 671 Lee, J. U., Schwan, E., & Meyer, C. M.
672 2019. Manipulating the Difficulty of C-Tests.
673 In *Proceedings of the 57th Annual Meeting of the*
674 *Association for Computational Linguistics*, pages
675 360–370.
- 676 Lee, J. U., Meyer, C. M., & Gurevych, I. 2020.
677 Empowering active learning to jointly optimize
678 system and user demands. *arXiv preprint*
679 *arXiv:2005.04470*.
- 680 Levenshtein, V. I. 1966. Binary codes capable of
681 correcting deletions, insertions, and reversals.
682 In *Soviet physics doklady* Vol. 10, No. 8, pages 707-
683 710.
- 684 Matsumori, S., Okuoka, K., Shibata, R., Inoue, M.,
685 Fukuchi, Y., & Imai, M. 2023. Mask and Cloze:
686 Automatic Open Cloze Question Generation Using
687 a Masked Language Model. *IEEE Access*, 11, pages
688 9835-9850.
- 689 Panda, S., Gomez, F. P., Flor, M., & Rozovskaya, A.
690 2022. Automatic Generation of Distractors for Fill-
691 in-the-Blank Exercises with Round-Trip Neural
692 Machine Translation. In *Proceedings of the 60th*
693 *Annual Meeting of the Association for*
694 *Computational Linguistics: Student Research*
695 *Workshop*, pages 391–401.
- 696 Qiu, Z., Wu, X., & Fan, W. 2020. Automatic distractor
697 generation for multiple choice questions in standard
698 tests. *arXiv preprint arXiv:2011.13100*.
- 699 Ren, S., & Zhu, K. Q. 2021. Knowledge-driven
700 distractor generation for cloze-style multiple choice
701 questions. In *Proceedings of the AAAI conference on*
702 *artificial intelligence*, Vol. 35, No. 5, pages 4339-
703 4347
- 704 Sakaguchi, K., Arase, Y., & Komachi, M. (2013,
705 August). Discriminative approach to fill-in-the-
706 blank quiz generation for language learners.
707 In *Proceedings of the 51st Annual Meeting of the*
708 *Association for Computational Linguistics (Volume*
709 *2: Short Papers)* (pp. 238-242).
- 710 Sireci, S.G. (1992). The Utility of IRT in Small-Sample
711 Testing Applications.
- 712 Smith, S.; Avinesh, P.; and Kilgarriff, A. 2010. Gap-fill
713 tests for language learners: Corpus-driven item
714 generation. In *Proceedings of ICON-2010: 8th*
715 *International Conference on Natural Language*
716 *Processing*, pages 1–6.
- 717 Susanti, Y. Iida, R. and Tokunaga, T. 2015. Automatic
718 Generation of English Vocabulary Tests.
719 In *Proceedings of the 7th International Conference*
720 *on Computer Supported Education*, pages 77-87.
- 721 Susanti, Y., Tokunaga, T., Nishikawa, H., & Obari, H.
722 2017. Controlling item difficulty for automatic
723 vocabulary question generation. *Research and*
724 *practice in technology enhanced learning*, 12(1):1–
725 16.
- 726 Uto, M., Tomikawa, Y., & Suzuki, A. (2023, July).
727 Difficulty-controllable neural question generation
728 for reading comprehension using item response
729 theory. In *Proceedings of the 18th Workshop on*
730 *Innovative Use of NLP for Building Educational*
731 *Applications (BEA 2023)* (pp. 119-129).
- 732 Wang, H. J., Hsieh, K. Y., Yu, H. C., Tsou, J. C., Shih,
733 Y. A., Huang, C. H., & Fan, Y. C. 2023. Distractor
734 Generation based on Text2Text Language Models
735 with Pseudo Kullback-Leibler Divergence
736 Regulation. In *Findings of the Association for*
737 *Computational Linguistics: ACL 2023*, pages
738 12477–12491.
- 739 Wojatzki, M., Melamud, O., & Zesch, T.
740 2016. Bundled Gap Filling: A New Paradigm for
741 Unambiguous Cloze Exercises. In *Proceedings of*
742 *the 11th Workshop on Innovative Use of NLP for*
743 *Building Educational Applications*, pages 172–181.
- 744 Xie, Q., Lai, G., Dai, Z., & Hovy, E. 2018. Large-scale
745 Cloze Test Dataset Created by Teachers.
746 In *Proceedings of the 2018 Conference on*
747 *Empirical Methods in Natural Language*
748 *Processing*, pages 2344–2356.
- 749 Yeung, C. Y., Lee, J. S., & Tsou, B. K. 2019. Difficulty-
750 aware Distractor Generation for Gap-Fill Items. In
751 *Proceedings of the The 17th Annual Workshop of the*
752 *Australasian Language Technology Association*,
753 pages 159–164.
- 754 Zaheer, M., Guruganesh, G., Dubey, K. A., Ainslie, J.,
755 Alberti, C., Ontanon, S., ... & Ahmed, A. 2020. Big
756 bird: Transformers for longer sequences. *Advances*
757 *in Neural Information Processing Systems*, 33,
758 pages 17283-17297.
- 759 Zesch, T., & Melamud, O. 2014. Automatic Generation
760 of Challenging Distractors Using Context-Sensitive
761 Inference Rules. In *Proceedings of the Ninth*
762 *Workshop on Innovative Use of NLP for Building*
763 *Educational Applications*, pages 143–148.

770 A Data Statistics

771 Table 7 presents the number of items per split in
772 our dataset.

773

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

774

775 Table 7: Data Statistics

776

776 B Algorithms

777 Figures 8, 9 present the algorithms Confidence-
778 Ranking Control, and 3-Factor Ranking Control
779 respectively.

Algorithm 1: Distractor Generation with Confidence Ranking

Input : A sentence S with a cloze, Answer Word A , Pretrained Model M , Proficiency Vocabulary List V , Original Options $Options$

Output: Difficult distractors D_{hard} , Easy distractors D_{easy}

```

1  $D_{hard}, D_{easy} \leftarrow \{\}, \{\}$ ;
2  $CandidateList \leftarrow sorted(ModelPredict(S, M, V))$ ; // Sort
   proficiency vocabulary  $V$  based on the scores predicted by
   the PTM  $M$ 
3  $OptionsRankList \leftarrow []$ ;
4 for  $opt$  in  $Options$  do
5    $OptionsRankList.append(CandidateList.index(opt))$ ;
6 end
7  $Rank_{hard} \leftarrow CandidateList.index(A)$ ;
8  $Rank_{easy} \leftarrow min(OptionsRankList)$ ;
9  $D_{hard} \leftarrow CandidateList[Rank_{hard}+1: Rank_{hard}+4]$ ;
10  $D_{easy} \leftarrow CandidateList[Rank_{easy}+1: Rank_{easy}+4]$ ;
11 return  $D_{hard}, D_{easy}$ 

```

780

781 Figure 8: Distractor Generation with Confidence-Ranking Control

782

Algorithm 2: Distractor Generation with 3-Factor Ranking

Input : A sentence S with a cloze, Answer Word A , Pretrained Model M , Proficiency Vocabulary List V , Original Options $Options$, Numbers of Candidate K

Output: Difficult distractors D_{hard} , Easy distractors D_{easy}

```

1  $D_{hard}, D_{easy} \leftarrow \{\}, \{\}$ ;
2  $CandidateList \leftarrow sorted(ModelPredict(S, M, V))$ ; // Sort
   proficiency vocabulary  $V$  based on the scores predicted by
   the PTM  $M$ 
3  $GloveSimList, LevenSimList \leftarrow [], []$ ;
4 for  $word$  in  $CandidateList$  do
5    $G \leftarrow Calculate\ Glove\ similarity(A, word)$ ;
6    $L \leftarrow Calculate\ Leven\ similarity(A, word)$ ;
7    $GloveSimList.append(word, G)$ ;
8    $LevenSimList.append(word, L)$ ;
9 end
10  $Rank_{hard} \leftarrow CandidateList.index(A)$ ;
11  $OptionsRankList \leftarrow []$ ;
12 for  $opt$  in  $Options$  do
13    $OptionsRankList.append(CandidateList.index(opt))$ ;
14 end
15  $Rank_{easy} \leftarrow min(OptionsRankList)$ ;
16  $GloveSimList_{hard}, GloveSimList_{easy} \leftarrow GloveSimList[Rank_{hard}+1:
   Rank_{hard}+K], GloveSimList[Rank_{easy}+1: Rank_{easy}+K]$ ;
17  $LevenSimList_{hard}, LevenSimList_{easy} \leftarrow LevenSimList[Rank_{hard}+1:
   Rank_{hard}+K], LevenSimList[Rank_{easy}+1: Rank_{easy}+K]$ ;
18  $D_{hard} \leftarrow Top\ 2\ candidates\ by\ GloveSimList_{hard}, Top\ 1\ candidate\ by
   LevenSimList_{hard}$ ;
19  $D_{easy} \leftarrow Top\ 2\ candidates\ by\ GloveSimList_{easy}, Top\ 1\ candidate\ by
   LevenSimList_{easy}$ ;
20 return  $D_{hard}, D_{easy}$ 

```

783

784 Figure 9: Distractor Generation with 3-Factor Ranking Control

785 C Annotation for Invalid Distractor Control

786

787 We analyzed the issues of invalid distractors with
788 human evaluation. We recruited 9 college students
789 at the CET-6 English proficiency level as
790 annotators. The annotators work with our research
791 lab on regular basis and receive subsidy for their
792 annotation work under supervision of our
793 administrative office.

794 The invalid distractors will most likely appear
795 when generating hard items. Using BERT’s
796 confidence score ranking without validity control,
797 we generated distractors for 4,575 items randomly
798 selected from the CLOTH-H dataset. Manual
799 annotation identified 1,676 items as having at
800 least one invalid generated distractor (i.e., a
801 distractor that could fit as an answer in the gap).
802 As our control strategies involves ranking
803 distractors after the answer, we identified 302
804 items to further test validity rules. Among the 906
805 distractors generated, 482 were annotated as
806 invalid, representing an invalidity ratio of 53.2%.
807 After applying the Confidence-Ranking Control
808 method and 3-Factor Ranking Control method,
809 the ratios dropped to 20.3% and 17.3%
810 respectively (Table 8).

811

812

813

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%

814 Table 8. Manual annotation of 906 distractors
815 generated with confidence ranking w/o validity rules,
816 and our methods of Confidence-Ranking Control and
817 3-Factor Ranking Control

818 The following are examples of items with the
819 answer (bolded) and invalid distractors
820 (italicized) generated by confidence ranking
821 without validity rules. The same item with
822 distractors generated using Confidence-ranking
823 Control and 3-Factor Ranking Control is also
824 shown below:

825 **Example #1:**

826 I hope I did the right thing, Mom, Alice said. I saw
827 a cat, all bloody but alive. I [MASK] it to the vet's,
828 and was asked to make payment immediately.

829 (1) Original options:

830 A. **carried** B. followed C. returned D. guided

831 (2) Distractors generated without control:

832 A. **carried** B. *took* C. brought D. delivered

833 (3) Distractors generated with 3-Factor Ranking
834 Control:

835 A. **carried** B. showed C. reported D. tried

836 (4) Distractors generated with Confidence

837 Ranking Control:

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

839 **Example #2:**

840 [MASK] this surprised him very much, he went
841 through the paper twice, but was still not able to
842 find more than one mistake, so he sent for the
843 student to question him about his work after the
844 exam.

845 (1) Original options:

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

847 (2) Distractors generated without control:

848 A. **As** B. *Because* C. Although D. Though

849 (3) Distractors generated with 3-Factor Ranking
850 Control:

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

852 (4) Distractors generated with Confidence
853 Ranking Control:

854 A. **As** B. Realizing C. Again D. Initially

855
856 **D Instruction to Annotators for Invalid**
857 **Distractor Identification.**

858 **Instruction:** You are given a set of multiple-
859 choice cloze test questions, each with four options.
860 The correct answer is identified, along with three
861 generated distractor options. Please review the
862 choices and identify any "invalid distractors" -
863 alternatives that contextually fit the gap as a
864 potentially correct response, rather than an
865 implausible one.

866 For example:

867 -----

868 When I began planning to move to Auckland to
869 study, my mother was worried about a lack of jobs
870 and cultural differences. Ignoring these ____ I got
871 there in July 2010.

872 A. **concerns** B. *worries*

873 C. fears D. considerations

874 -----

875 Here, the answer is "concerns". The generated
876 distractors include "worries". Both are
877 grammatically correct. "Concerns" fits the
878 semantic context only slightly better. Therefore,
879 in this case, "worries" is considered an "invalid
880 distractor".

881 Your annotation results will help assess the
882 efficacy of our difficulty-control strategies in
883 limiting invalid distractor generation for multiple
884 choice cloze tests.

885
886 **E Instruction to Annotators for Invalid**
887 **Item Difficulty Comparisons.**

888
889 **Instruction:** You are given a set of multiple-
890 choice cloze test questions, each with two sets of
891 options. Please compare and mark the harder set
892 with 2 and easy set with 1. Your annotation results
893 will help assess the efficacy of our difficulty-
894 control strategies for cloze item generation.

895

896

897 For example:

898 I have a good friend at school. Her name is Liu
899 Mei. She's fifteen years old. She is a beautiful girl
900 [1] bright eyes and long black hair. In some ways
901 we look the same, [2] some students say we are
902 twins. In our class, her math is not good. But she
903 works hard all the time. Now she is doing [3] than
904 before. I hope she can make great progress. I often
905 go to her house. There are many kinds of books
906 and magazines on her bookshelf. She likes
907 reading. Her Chinese is best in our class. She often
908 helps me with Chinese. Liu Mei is an active girl.
909 She's a little more outgoing than me. She likes
910 tennis very much. She is well at tennis. But I'm [4]
911 better than her at ping-pong.

Your score		Your score	
	['and', 'with', 'have', 'has']		['of', 'with', 'sporting', 'having']
	['or', 'so', 'until', 'including']		['because', 'so', 'or', 'but']
	['good', 'well', 'better', 'best']		['worse', 'harder', 'better', 'poorer']
	['zero', 'more', 'little', 'although']		['much', 'more', 'many', 'lot']

912