

Improving Automatic Grammatical Error Annotation for Chinese Through Linguistically-Informed Error Typology

Yang Gu¹ Zihao Huang¹ Min Zeng²
Mengyang Qiu^{1,3} Jungyeul Park^{1,4}

¹Open Writing Evaluation, France ²TEKsystems, Canada

³Department of Psychology, Trent University, Canada

⁴Department of Linguistics, The University of British Columbia, Canada

<http://open-writing-evaluation.github.io>

Introduction

- What is GEC?
 - Detects and corrects grammatical errors in text
 - Crucial for language learning and writing improvement
- Why focus on Chinese?
 - Growing demand for high-quality Chinese writing support tools
 - Chinese poses unique challenges due to its writing system
- Existing annotation tools:
 - ERRANT (English): Provides detailed annotations (Bryant et al., 2017)
 - ChERRANT: Adapted to Chinese but limited to basic character-level errors (Zhang et al., 2022)
- What this study does:
 - Refines the ChERRANT framework
 - Introduces a linguistically-informed typology for Chinese-specific errors for both L1 and L2 learners

Features of the Chinese Writing System

- No explicit word boundaries (e.g., continuous text without spaces)
 - Words can be 1 – 2+ characters
 - Example: 被 (*bèi*) can be:
 - Standalone passive marker
 - Part of words like 被子 (*bèizi*, “blanket”)
- Many homophones (e.g., 权力 vs. 权利) and visually similar characters (e.g., 四 vs. 西)
- Implications for GEC:
 - Tokenization errors
 - Misidentification of grammatical roles and structures

Previous CGEC Annotation System–ChEERRANT

- Examined:
 - Manual vs. automated annotations
 - L1 vs. L2 writing via automated annotations
- Validation sets of FCGEC (Xu et al., 2022) and MuCGEC (Zhang et al., 2022)

Dataset	# of sentences		Different numbers of errors per sentence					
	Ori.	Ref.	0	1	2	3	4	≥ 5
FCGEC (L1)	2000	2143	899	1138	101	3	2	0
MuCGEC (L2)	1138	1153	55	214	238	176	153	317

Previous CGEC Annotation System–ChERRANT

- Manual vs. automated annotations
 - Automated annotations align with manual annotations in up to 95% of cases, with a match rate of 73% for more complex sentences containing multiple grammatical errors
 - Differences in word boundary may lead to different error type definitions over different spans of text
 - Problems of dealing with changing the order of large spans of text like phrases and clauses
- L1 vs. L2 writing
 - The ChERRANT analysis results are consistent with conventional assumptions about the differences between L1 and L2 writing
 - ChERRANT is less reliable when handling sentences with multiple complex errors (e.g., a potential underrepresentation of ordering errors)

Refined Annotation for CGEC

- A new error typology that better captures both L1 and L2 grammatical errors, along with a new implementation based on this scheme
 - Similar pronunciation
 - Similar shapes
 - Multifaceted similarity
 - The *de* particles—的, 地, and 得
 - Character order

Refined Annotation for CGEC

Learner writing	Correction
(1) 每个抽烟的人都有这样的 权力 。	每个抽烟的人都有这样的 权利 。
'Every smoker has this power .'	'Every smoker has this right .'
(2) 交朋友的 时后 ，很可能会碰到矛盾。	交朋友的 时候 ，很可能会碰到矛盾。
'Time after making friends, it's highly likely to run into conflicts.'	'When making friends, it's highly likely to run into conflicts.'
(3) 我 一前 没住过五星级旅馆，所以我很惊讶。	我 以前 没住过五星级旅馆，所以我很惊讶。
'I haven't stayed in five-star hotels one before , so I am very surprised.'	'I haven't stayed in five-star hotels before , so I am very.'
(4) 她有两个姐姐、一个妹妹和 西 个哥哥。	她有两个姐姐、一个妹妹和 四 个哥哥。
'She has two older sisters, one younger sister, and west older brothers.'	'She has two older sisters, one younger sister, and four older brothers.'
从十六世纪开始， 欧州人 就抽烟。	从十六世纪开始， 欧洲人 就抽烟。
'Since the 16th century, Europe region people smoked.'	'Since the 16th century, Europeans smoked.'
(6) 经理， 新得 计划发您信箱了，您看了吧？	经理， 新的 计划发您信箱了，您看了吗？
'Manager, the new obtain plan was sent to your mailbox, you've seen it, right?'	'Manager, the new MOD plan was sent to your mailbox, have you seen it?'
(7) 简单生活，哪怕对身体还是精神，还大有裨益。	简单 的 生活，无论对身体还是精神，都大有裨益。
'A simple life, whether for the body or the mind, is greatly beneficial.'	'A simple MOD life, whether for the body or the mind, is greatly beneficial.'
(8) 反而那些不帅，还有点丑但是很会唱歌就被淘汰了。	反而那些不帅，还有点丑但是很会唱歌 的 就被淘汰了。
'Instead, those who are not handsome, a bit ugly, but can sing well were eliminated.'	'Instead, those who are not handsome, a bit ugly, but can sing well were eliminated.'
(9) 要了解一个人，不妨看他读些 么什 书，观察向他来往得朋友一样有效。	要了解一个人，不妨看他读些 什么 书，这跟观察与他来往的朋友一样有效。
'To understand a person, it's just as effective to see what books he reads and observe the friends he interacts with.'	'To understand a person, it's just as effective to see what books he reads and observe the friends he interacts with.'

Refined Annotation for CGEC

- Stanza for segmentation and tagging

Algorithm 1 Pseudo-code for error classification

```
1: function ERRORCLASSIFICATION ( $\mathcal{S}, \mathcal{T}, \{R|M|U\}$ ):
2:   if  $R \wedge (\text{SIM}(\text{pinyin}) > \alpha_1) \wedge (\text{SIM}(\text{shape}) > \alpha_2)$ 
   then
3:     return R:MULTI
4:   else if  $R \wedge (\text{SIM}(\text{pinyin}) > \alpha_1)$  then
5:     if  $\mathcal{T} == \text{de}$  then
6:       return R:DE
7:     else
8:       return R:PINYIN
9:     end if
10:  else if  $R \wedge (\text{SIM}(\text{shape}) > \alpha_2)$  then
11:    return R:SHAPE
12:  else if  $R \wedge (\text{SET}(\mathcal{S}) == \text{SET}(\mathcal{T}))$  then
13:    if  $\text{LEN}(\mathcal{T}) == 1$  then
14:      return R:CO
15:    else
16:      return R:WO
17:    end if
18:  else if  $M \wedge (\mathcal{T} == \text{de})$  then
19:    return M:DE
20:  end if
21:  return  $\{R|M|U\}$ 
```

Refined Annotation for CGEC

- Identified 12-16% more errors than ChERRANT
- Greater detail in annotating spelling and structural errors

Error Type	L1		L2	
	Count	Ratio	Count	Ratio
M:DE	6	0.0051	8	0.0207
R:DE	5	0.0032	14	0.0031
R:MULTI	73	0.0466	197	0.0433
R:PINYIN	21	0.0134	146	0.0321
R:SHAPE	20	0.0128	164	0.0361
R:WO	63	0.0402	148	0.0326
Total	188	0.1212	677	0.1678

Conclusion and Future Work

- Key contributions:
 - Introduced a refined error typology tailored to Chinese
 - Highlighted differences in L1 and L2 errors
- Future directions:
 - Address more context-sensitive grammatical issues in addition to surface-level errors
 - Improve robustness of segmentation and annotation tools

Thank you!



References

- Bryant, C., Felice, M., and Briscoe, T. (2017). Automatic Annotation and Evaluation of Error Types for Grammatical Error Correction. In Barzilay, R. and Kan, M.-Y., editors, *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 793–805, Vancouver, Canada. Association for Computational Linguistics.
- Xu, L., Wu, J., Peng, J., Fu, J., and Cai, M. (2022). FCGEC: Fine-Grained Corpus for Chinese Grammatical Error Correction. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1900–1918, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Zhang, Y., Li, Z., Bao, Z., Li, J., Zhang, B., Li, C., Huang, F., and Zhang, M. (2022). MuCGEC: a Multi-Reference Multi-Source Evaluation Dataset for Chinese Grammatical Error Correction. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3118–3130, Seattle, United States. Association for Computational Linguistics.