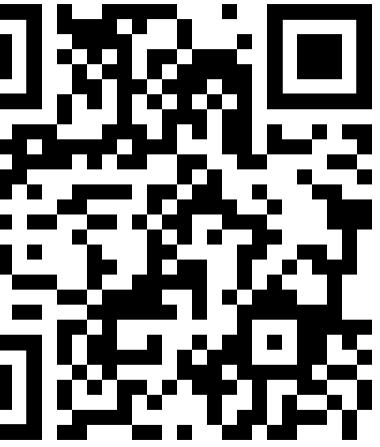
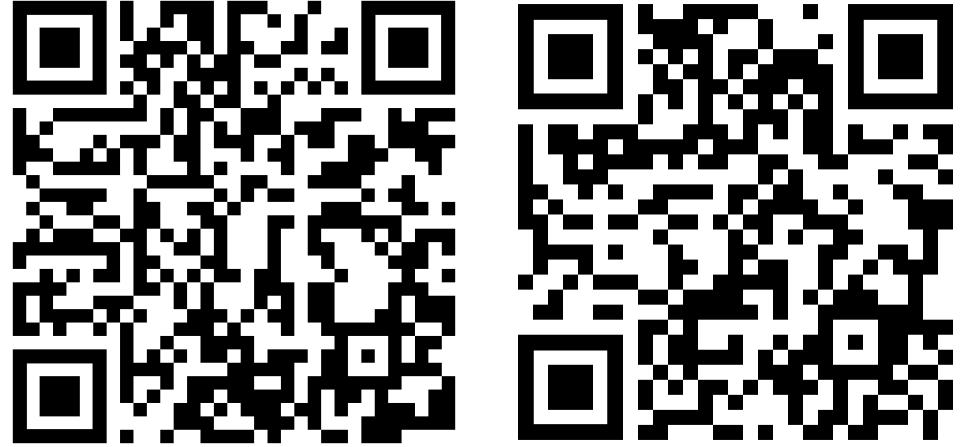


Towards Standardizing Korean Grammatical Error Correction: Datasets and Annotation

Soyoung Yoon¹, Sungjoon Park², Gyuwan Kim³, Junhee Cho⁴, Kihyo Park⁵, Gyutae Kim², Minjoon Seo¹, Alice Oh⁶

KAIST AI¹, Softly AI², University of California, Santa Barbara³, Cornell University⁵, Google⁴, KAIST⁶



Introduction

Motivation: Research on Korean Grammatical Error Correction(GEC) is limited. We attribute this problem to the lack of carefully designed evaluation benchmark and resource for Korean GEC.

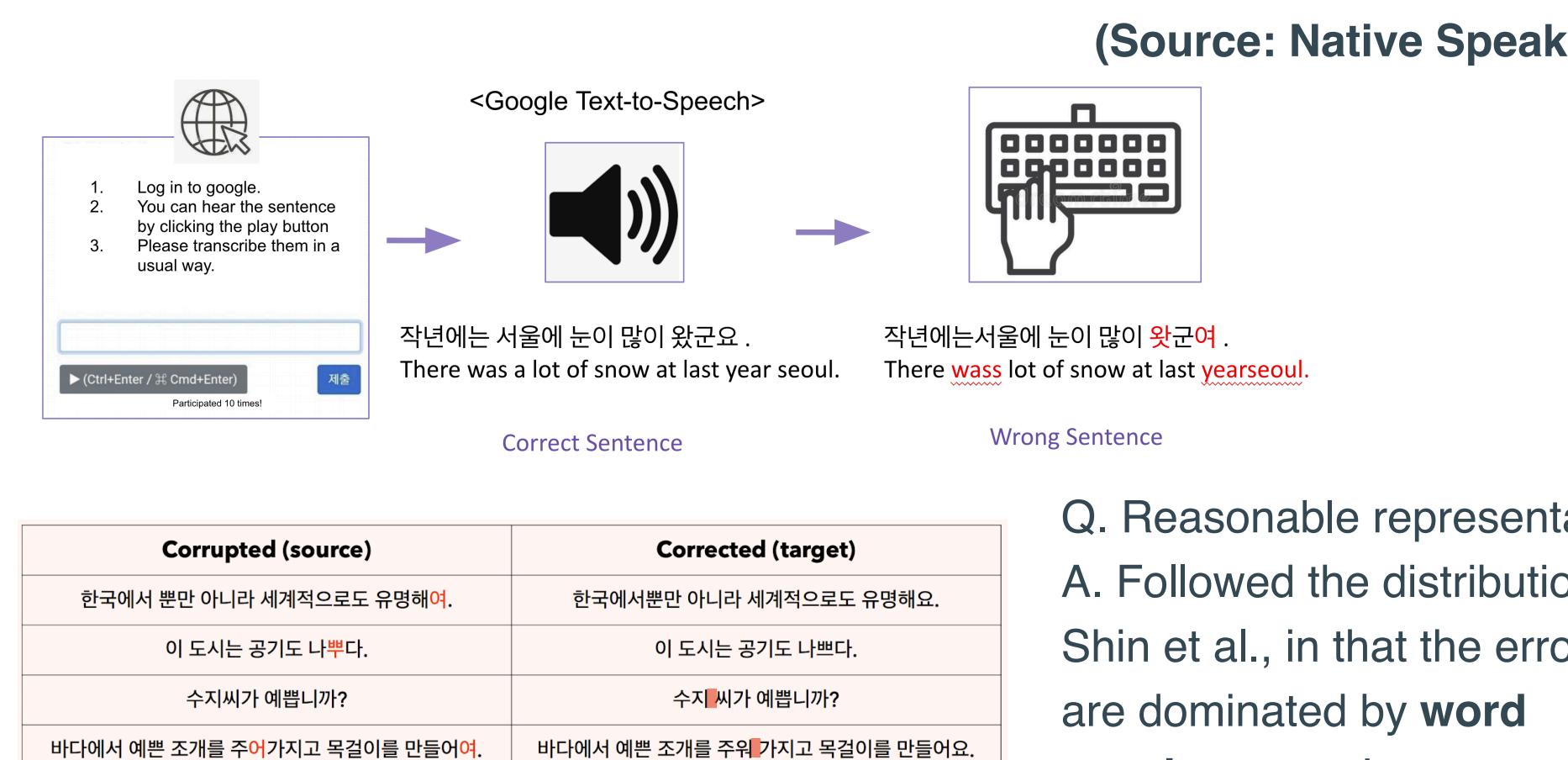
Contribution:

- Three large-scale organized **dataset** of Korean GEC
- KAGAS, an **automatic annotation & evaluation system** for Korean GEC
- Open-sourced baseline **error correction model** based on our data

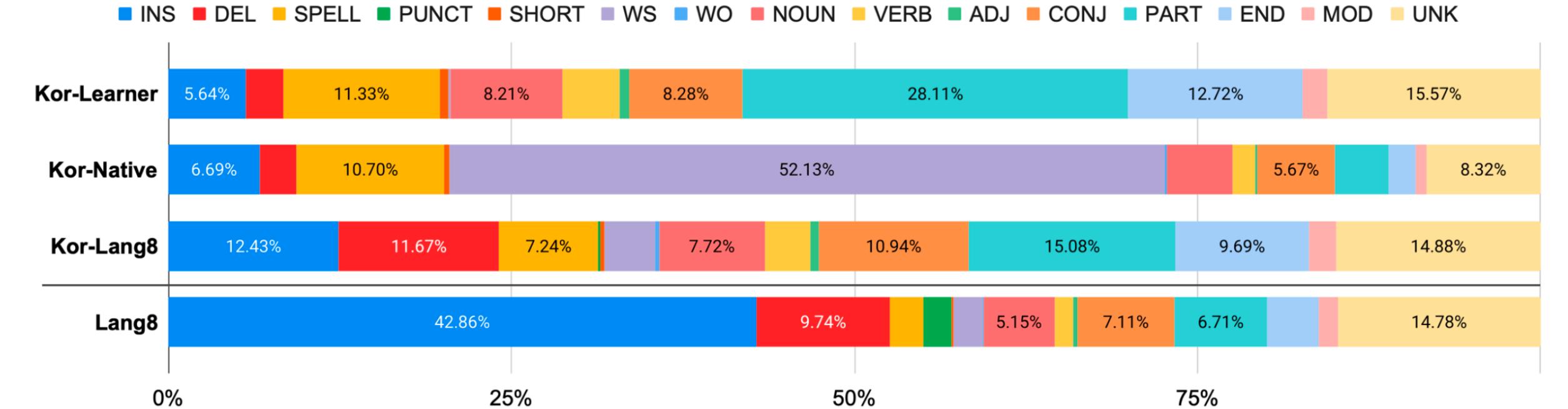
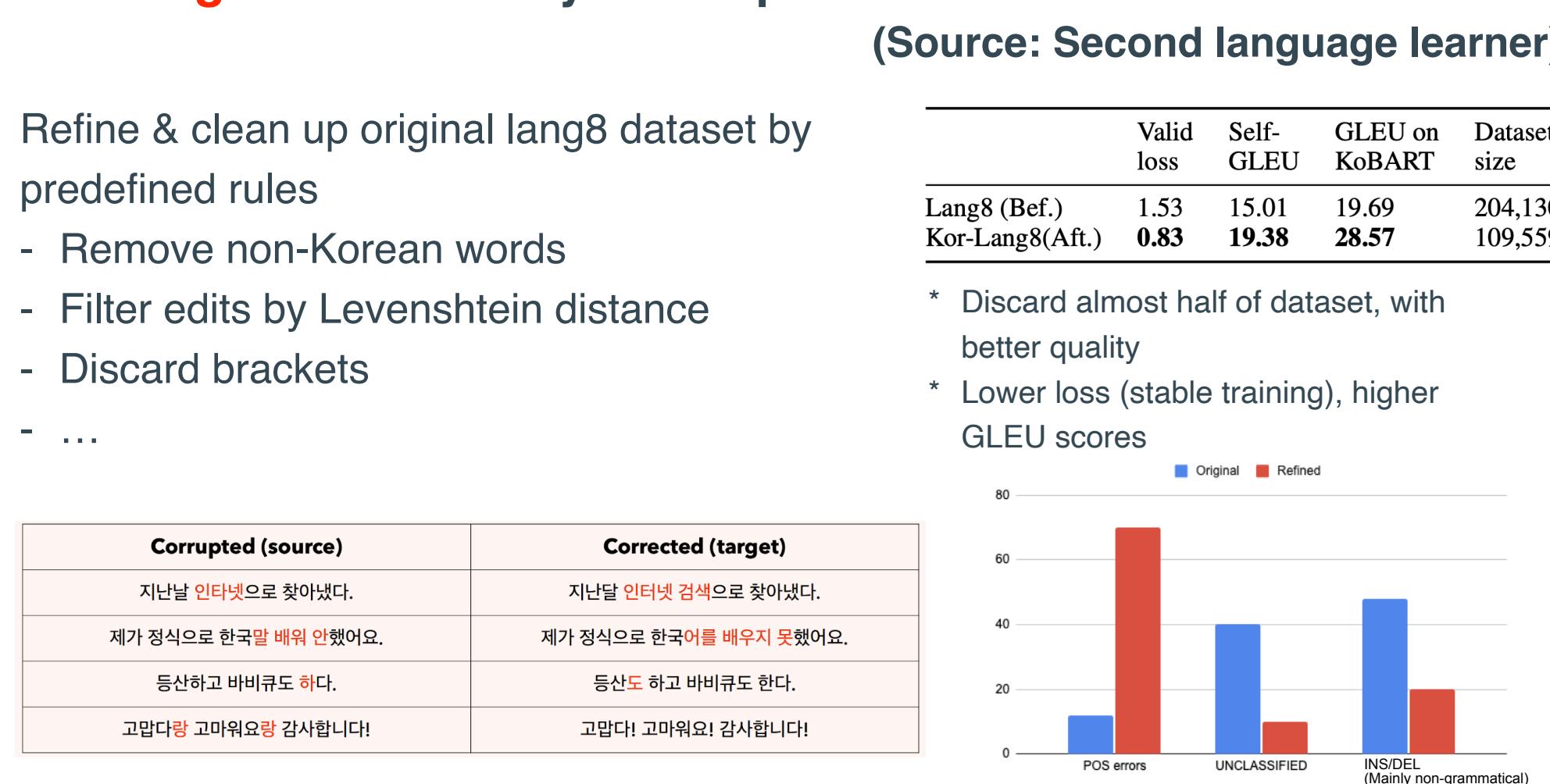
Dataset

	Kor-Learner	Kor-Native	Kor-Lang8
# Sentence pairs	28,426	17,559	109,559
Avg. token length	14.86	15.22	13.07
# Edits	59,419	29,975	262,833
# Edits / sentence	2.09	1.71	2.40
Avg. tokens per edit	0.97	1.40	0.92
Prop. tokens changed	28.01%	29.37%	39.42%

Kor-Native - data collection by transcription



Kor-Lang8 - collected by social platform



Kor-Learner - essays & annotations by Korean learners and tutors

Merge morpheme-level edits into word-level edits to make parallel corpora

Apply Korean orthography guidelines



(Source: Second language learner)

Kor-Learner contains high proportion of particle errors

Corrupted (source)	Corrected (target)
우리 모두 꿈을 잘 이루었으면 좋겠습니다.	우리 모두의 꿈이 잘 이루어졌으면 좋겠습니다.
그 제 꿈이 교수님은 것입니다.	제 꿈은 교수가 되는 것입니다.
밥을 먹으려 식당에 가요.	밥을 먹으려 식당에 가요.
그래서 집 관리비를 절약할 수 있어요.	그래서 집 관리비를 절약할 수 있어요.

* original XML file:

```
<LearnerErrorAnnotations>
<word>
<w>과 드 데 인</w>
<morph from="178" subsequence="1" to="182" wordStart="Start">
<Proread pos="NNG" />
<ErrorArea pos="CN" />
<ErrorPattern type="REP" />
<ErrorLevel type="MDV/POS" />
</morph>
<morph from="178" subsequence="2" to="182" wordStart="None">
<Proread pos="XSV" />
<ErrorArea pos="CN" />
<ErrorPattern type="MDV/POS" />
<ErrorLevel type="MDV/POS" />
</morph>
<morph from="178" subsequence="4" to="182" wordStart="None">
<Preserve />
<Proread pos="XSV" />
<ErrorArea pos="CN" />
<ErrorPattern type="MDV/POS" />
<ErrorLevel type="MDV/POS" />
</morph>
</word>
</LearnerErrorAnnotations>
```

과도하ㄴ -> 과도한

KAGAS

We build an **automatic annotation toolkit** that annotates parallel Korean sentences with error type information given **original & corrected pairs** considering Korean linguistics

Motivation: Annotation by humans have disadvantages

What is different for Korean?

- Classifying error types by word morpheme-level

- Morpheme-level deletion
 - 소풍(NNG) + 을(JKO) => 소풍(NNG)
 - 을(JKO) is deleted, thus labeled as PART(JKO is grouped to PART)
- Morpheme-level insertion
 - 유학(NNG) => 유학(NNG) + 러(NNP)
 - 러(NNP) inserted, thus labeled as Noun(NNP -> PART)

Q. Reasonable representative?

A. Followed the distribution of Shin et al., in that the errors are dominated by **word spacing** errors!

- No PREP, but PART for Korean

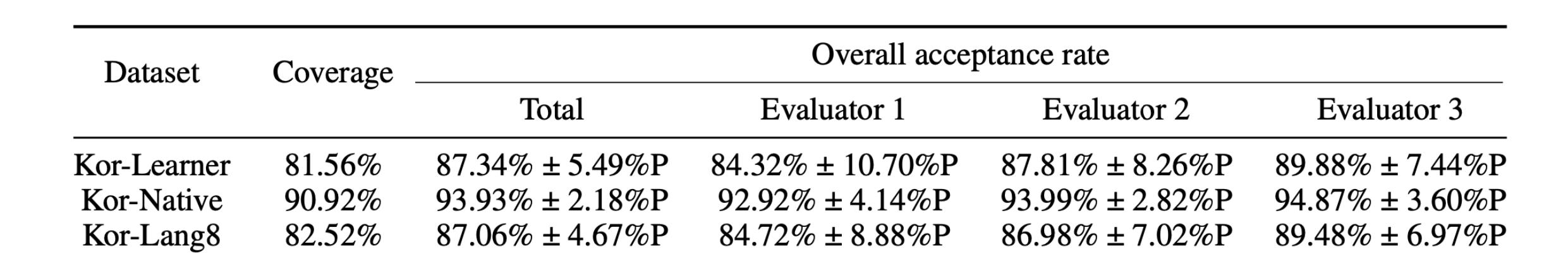
School->학교, To->에, 학교에->School-To(not To-School)

"To" is not preposition, but rather a postpositional **particle**

- About INS/DEL edits.

Korean is a discourse-oriented language. One can omit the subject or object in a sentence depending on the previous context, which is grammatically correct in most cases.

Acceptance rate of KAGAS by 3 GEC experts



Error Code	Description & Acceptance Rate (%)	Example
INS	A word is inserted. 100.00% ± 0.00%P	고등학교 때 어긴 경험 Original: 고등학교 때 어긴 경험 Corrected: 고등학교 때 규칙을 어긴 경험 Translation: Experience to break a rule in high school
DEL	A word is deleted. 100.00% ± 0.00%P	전쟁 끝 직후 장군들은 사형을 선고 받았다. Original: 전쟁 끝 직후 장군들은 사형을 선고 받았다. Corrected: 전쟁 직후 장군들은 사형을 선고 받았다. Translation: After the war, the generals are sentenced to death.
WS	Spacing between words is changed. 100.00% ± 0.00%P	이 옷은 더러워요. Original: 이 옷은 더러워요. Corrected: 이 옷은 더러워요. Translation: This cloth is dirty.
WO	The order of words is changed. 97.44% ± 3.51%P	저는 더 한국어를 배우고 싶어요. Original: 저는 더 한국어를 배우고 싶어요. Corrected: 저는 한국어를 더 배우고 싶어요. Translation: I want to learn Korean further.
SPELL	Spelling error 97.44% ± 3.51%P	파티에서 우리는 춤을 췄요. Original: 파티에서 우리는 춤을 췄요. Corrected: 파티에서 우리는 춤을 췄요. Translation: We dance at the party.
PUNCT	Punctuation error 98.72% ± 2.50%P	1993년 의 겨울의 일이었다. Original: 1993년 의 겨울의 일이었다. Corrected: 1993년, 겨울의 일이었다. Translation: It was 1993, a happening in winter.
SHORT	An edit that does not change the structure of morphemes. 73.08% ± 9.84%P	한국어는 저한테 너무 어려운 언어이었어요. Original: 한국어는 저한테 너무 어려운 언어이었어요. Corrected: 한국어는 저한테 너무 어려운 언어였어요. Translation: Korean Language was too difficult to me.
VERB	An error on verb 79.49% ± 8.96%P	어제 친구에게 편지를 썼어요. Original: 어제 친구에게 편지를 썼어요. Corrected: 어제 친구에게 편지를 썼어요. Translation: I wrote a letter to my friend yesterday.
ADJ	An error on adjective 73.08% ± 9.84%P	친한 친구 친한 친구 Original: 친한 친구 Corrected: 친한 친구 Translation: A close friend.
NOUN	An error on noun 75.64% ± 9.53%P	나중에 기회가 있을 때 한국에 유학하고 싶습니다. Original: 나중에 기회가 있을 때 한국에 유학하고 싶습니다. Corrected: 나중에 기회가 있을 때 한국에 유학하고 싶습니다. Translation: I want to study abroad in Korea in the future.
PART	An error on particle 97.44% ± 3.51%P	하와이에서 사는 우리 사촌 하와이에서 사는 우리 사촌 Original: 하와이에서 사는 우리 사촌 Corrected: 하와이에 사는 우리 사촌 Translation: My cousin living in Hawaii
END	An error on ending 87.18% ± 7.42%P	오래 기다려요. Original: 오래 기다렸어요. Corrected: 오래 기다렸어요. Translation: I waited for a long time.
MOD	An error on modifier 89.74% ± 6.73%P	점심이 너무 막은 나머지 배고팠어요. Original: 점심이 너무 막은 나머지 배고팠어요. Corrected: 점심이 너무 막은 나머지 배고팠어요. Translation: I was hungry because I had such a small lunch.
CONJ	An error on conjugation 43.59% ± 11.00%P	오늘은 머리를 잘라야 갔다. Original: 오늘은 머리를 자르러 갔다. Corrected: 오늘은 머리를 자르러 갔다. Translation: I went to a barber to get my hair cut today.

Full list of KAGAS Error Types with examples

Baseline Models

Trained our data on top of SKT-AI/KoBART and compared with Hanspell, a widely used GEC system based on statistical modeling.



한국어 맞춤법/문법 검사기

맞춤법/문법 검사하기

다시 쓰기

[종 글자 수]

Model variants	Kor-Learner			Kor-Native			Kor-Lang8			Kor-Union			Gen. time	
GLEU	M²	Pre.	Rec.	F_{0.5}	GLEU	M²	Pre.	Rec.	F_{0.5}	GLEU	M²	Pre.	Rec.	F_{0.5}

<tbl_r cells="14" ix="2" maxcspan="1" maxrspan="1