

# Evaluating Prompting Strategies for GEC Based on Language Proficiency

<https://github.com/jungyeul/prompting-gec> LREC-COLING 2024

## Prompting GEC

Min Zeng<sup>†\*</sup>, Jiexin Kuang<sup>†\*</sup>, Mengyang Qiu<sup>‡</sup>, Jayoung Song<sup>¶</sup>, Jungyeul Park<sup>‡</sup>

<sup>†</sup>The University of British Columbia, Canada, <sup>‡</sup>Department of Psychology, Trent University, Canada, <sup>¶</sup>Department of Asian Studies, Pennsylvania State University, USA. \*Min Zeng and Jiexin Kuang contributed equally.

### Results

	A						B						C						all					
	TP	FP	FN	Prec	Rec	F0.5	TP	FP	FN	Prec	Rec	F0.5	TP	FP	FN	Prec	Rec	F0.5	TP	FP	FN	Prec	Rec	F0.5
GPT-2 zero-shot	70	3944	2876	0.0174	0.0237	0.0184	45	5204	2453	0.0086	0.018	0.0096	28	4860	1058	0.0057	0.0258	0.0068	143	14008	6389	0.0101	0.0219	0.0113
GPT-3.5 zero-shot	1203	3770	1740	0.2419	0.4088	0.2634	940	4693	1556	0.1669	0.3766	0.1878	407	4183	677	0.0887	0.3755	0.1047	2550	12646	3973	0.1678	0.3909	0.1894
SOTA GECTOR T5	1046	632	2054	0.6234	0.3374	0.533	785	458	1836	0.6315	0.2995	0.5169	315	208	845	0.6023	0.2716	0.4843	2146	1298	4735	0.6231	0.3119	0.5194

Prompting results using GPT-2 (gpt2-x1 and FT = fine-tuned), GPT-3.5 (text-davinci-003) and SOTA results by models of GECTOR and T5.

### Analysis and Discussion

		TP	FP	FN	Prec	Rec	F0.5
M:PUNCT	A	189	171	134	0.525	0.5851	0.536
	B	203	132	133	0.606	0.6042	0.6056
	C	95	96	80	0.4974	0.5429	0.5059
R:VERB	A	21	60	113	0.2593	0.1567	0.2293
	B	17	55	113	0.2361	0.1308	0.2033
	C	6	43	51	0.1224	0.1053	0.1186
M	A	318	436	372	0.3703	0.3571	0.1691
	B	336	347	344	0.4919	0.4941	0.2458
	C	157	222	168	0.4142	0.4830	0.2180

#### 1. Label-by-label evaluation approach:

2. Is recall higher than precision in prompting GPT for the GEC task? Consistent higher recall compared to precision showcases a tendency of over-correction in prompting GPT for the GEC task. We have observed that proficiency levels A and B, however, do not exhibit such a propensity. It holds true even for GPT-3.5, where recall consistently surpasses precision. Nevertheless, the difference between precision and recall measurements in levels A and B is considerably smaller compared to level C.

#### 3. Results using various F-scores

	FT GPT-2			GPT-3.5		
	F0.5	F1	F2	F0.5	F1	F2
A	0.4192	0.4032	0.3885	0.3784	0.4030	0.4310
B	0.4210	0.4010	0.3827	0.3291	0.3625	0.4034
C	0.3310	0.3388	0.3470	0.2199	0.2680	0.3430
all	0.3907	0.4029	0.3792	0.3590	0.3230	0.4040

4. Comparison between prompting GPT and SOTA State-of-the-art (SOTA) results continue to demonstrate superior performance compared to prompting GPT in the GEC task in all aspects of results including precision and recall measures regardless of proficiency levels. Our assumption is primarily based on the fact that SOTA models are usually subjected to extensive fine-tuning processes.

Acknowledgement: This work was supported in part by Oracle Cloud credits and related resources provided by Oracle for Research.

### Experimental Settings:

Proficiency A		Proficiency B		Proficiency C	
M:PUNCT	0.0933	M:PUNCT	0.1134	M:PUNCT	0.1183
R:ORTH	0.0602	R:PREP	0.0589	R:PREP	0.0517
R:PREP	0.0506	M:DET	0.0442	M:DET	0.0345
R:VERB:TENSE	0.0455	R:VERB	0.0414	R:VERB	0.0323
R:VERB	0.0419	R:VERB:TENSE	0.0393	R:VERB:TENSE	0.0273

Most frequent errors and their ratio in W&I

model	gpt2-x1
tokenizer	gpt2-x1
num_exemplars	0-4 shots
max_model_token_length	256 if num_exemplars is 0 else 512
delimiter left and right	{ }

Experimental setting for GPT-2 (gpt2-x1) inferences, and we also adapt it to GPT-3.5 (text-davinci-003)

1-shot	ungrammatical	This is important thing.
2-shot	grammatical	This is an important thing.
2-shot	ungrammatical	Water is needed for alive.
3-shot	grammatical	Water is necessary to live.
3-shot	ungrammatical	And young people spend time more their lifestyle.
4-shot	grammatical	And young people spend more time on their lifestyles.
4-shot	ungrammatical	Both of these men have dealt with situations in an unconventional manner and the results are with everyone to see.
4-shot	grammatical	Both of these men have dealt with situations in an unconventional manner and the results are plain to see.

### Prompt examples

epochs	5
using masked language modeling	False
block size (train)	128
per_device_train_batch_size	4
save_steps	10000
save_total_limit	2

### Fine-tuning parameters

### Conclusion:

- We investigated the strengths and limitations of prompting GPT for the GEC task based on different language proficiency levels.
- We used our own implementations to calculate relevant metrics for label-by-label analysis.
- We observed a tendency of over-correction in prompting GPT, and it is more obvious in the recent version of GPTs, where recall consistently surpasses precision.



Download the paper →

