

Towards L2-friendly pipelines for learner corpora: A case of L2-Korean learners

SUMMARY

- We introduce the Korean-Learner-Morpheme (KLM) corpus, a manually annotated dataset from second language (L2) learners of Korean, featuring morpheme tokenization and part-of-speech (POS) tagging ($n = 129,784$).
- We evaluate the performance of four Korean morphological analyzers in tokenization and POS tagging on the L2-Korean corpus. Results highlight the analyzers' reduced performance on L2 data, indicating the limitation of advanced deep-learning models when dealing with L2-Korean corpora.
- Model fine-tuning with the KLM corpus demonstrates improved tokenization and POS tagging accuracy on the L2-Korean dataset.

DATASETS & ANNOTATION

L2 corpus: The KLM corpus, comprising 129,784 morphemes with morpheme tags grounded in the Sejong tag set

- The inclusion of data concerning classroom proficiency levels (ranging from 1 to 6 as a proxy for learner proficiency), nationality, gender, and writing topics
- The random extraction of 600 texts from the original corpus (Park & Lee, 2016), with each proficiency level represented by 100 texts
- The manual annotation of the corpus by three Korean native speakers, featuring detailed descriptions during Korean morpheme annotations and their evaluations

Category	Token	Tags
# of refinement	19,481	20,987
% of refinement	15.01	16.17
# of agreement	128,890	128,243
% of agreement	99.31	98.81
Total	129,784	

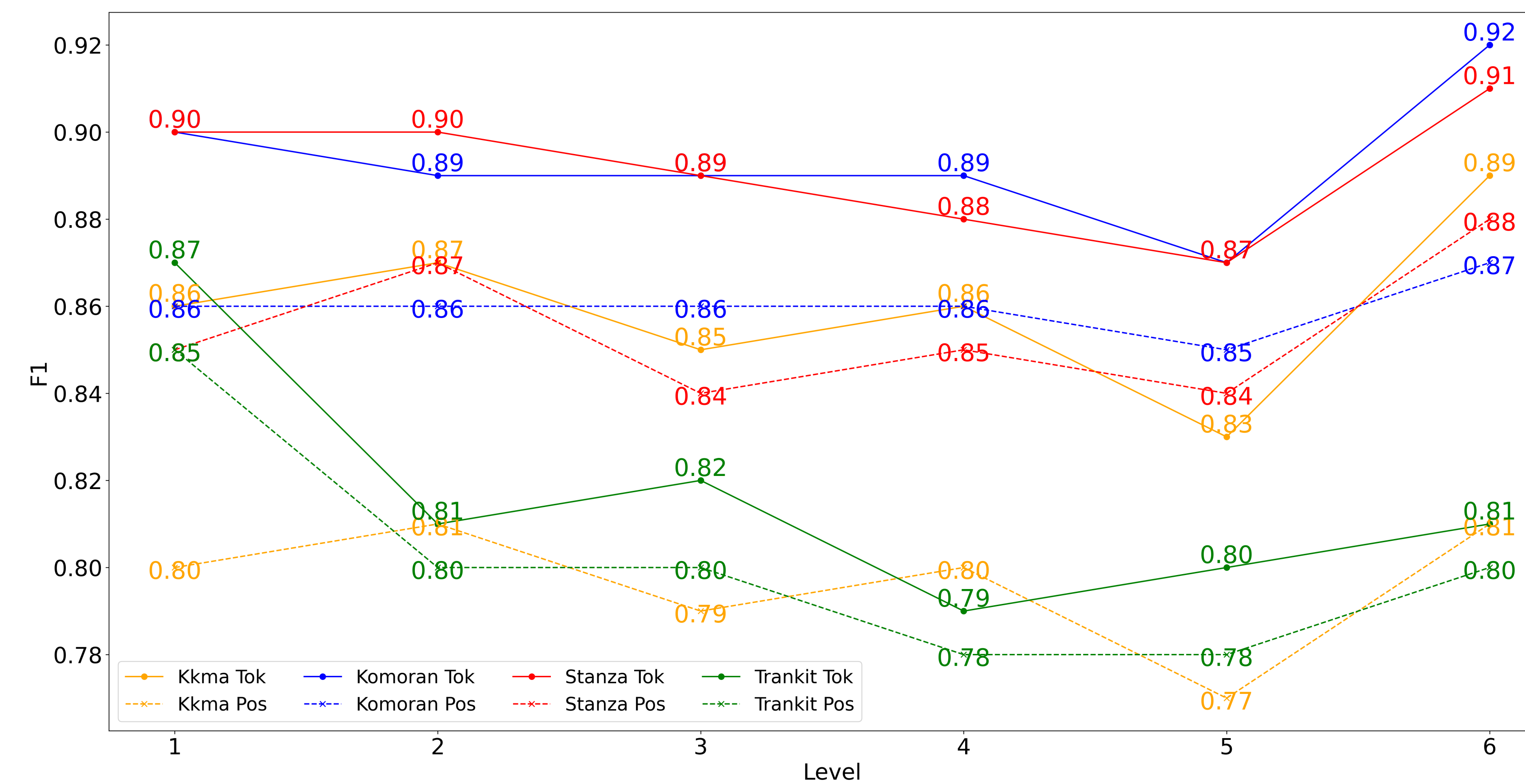
Reference L1 corpus: Google Korean Universal Dependency Treebank (UD Korean GSD) as a baseline

OVERALL PERFORMANCE

Morphological analyzers: The employment of four open-access morphological analyzers, incorporating various computational algorithms ranging from statistical models to deep-learning models (i.e., Stanza: Deep Biaffine attention for neural network; Trankit: Transformer; Kkma & Komoran: Hidden Markov model)

Analyzer	L2 TOK	L2 POS	L1 TOK	L1 POS
Stanza	0.89	0.86	0.92	0.93
Trankit	0.81	0.80	0.85	0.88
Kkma	0.86	0.80	0.88	0.81
Komoran	0.89	0.86	0.92	0.86

BY-LEVEL PERFORMANCE



MODEL TRAINING & RE-EVALUATION

- Training:** The Stanza model, further trained on an L2 dataset (KLM corpus), aiming to evaluate the potential improvement in performance compared to an exclusively L1-trained model
- Results:** The F1 scores for the Stanza+L2 model showing improvements with **tokenization** $0.93 > 0.89$ and **POS tagging** $0.91 > 0.86$ (compared to the highest scores of models trained exclusively on the L1 dataset)

BY-TAG PERFORMANCE

POS Tags	Analyzer		
	Stanza	Komoran	Stanza+L2
JKO	0.94 ⁽⁴⁷⁰⁵⁾	0.93 ⁽²²¹²⁾	0.96 ⁽⁴⁵⁴⁾
MAJ	0.94 ⁽¹¹⁹²⁾	0.94 ⁽⁶⁶⁸⁾	0.85 ⁽¹⁴³⁾
JKS	0.92 ⁽⁴¹⁶⁰⁾	0.91 ⁽¹⁸⁷⁴⁾	0.95 ⁽⁴⁰²⁾
JKG	0.92 ⁽¹²⁵⁷⁾	0.85 ⁽⁴²³⁾	0.95 ⁽¹¹⁹⁾
VCN	0.91 ⁽¹⁷⁸⁾	0.95 ⁽⁷⁵⁾	0.86 ⁽²⁶⁾
JKB	0.89 ⁽⁶³⁹⁹⁾	0.89 ⁽⁴²³⁾	0.92 ⁽⁶³⁴⁾
MAG	0.87 ⁽⁴⁶²⁸⁾	0.90 ⁽¹⁸⁸⁵⁾	0.86 ⁽⁴⁴⁶⁾
JX	0.86 ⁽⁵³¹⁷⁾	0.91 ⁽²³⁸⁴⁾	0.91 ⁽⁵⁴³⁾
NNB	0.85 ⁽⁴⁶⁸⁵⁾	0.84 ⁽¹⁸⁸⁷⁾	0.84 ⁽⁵³²⁾
XSN (Suffix, n.)	0.84 ⁽¹⁵⁵⁷⁾	0.85 ⁽⁵⁸¹⁾	0.87 ⁽¹³⁹⁾
ETN	0.83 ⁽⁸³¹⁾	0.89 ⁽³²⁶⁾	0.85 ⁽⁸³⁾
NNG (Noun, common)	0.77 ⁽³⁰³⁵³⁾	0.82 ⁽⁹⁶⁸²⁾	0.83 ⁽²⁸⁶⁶⁾
VCP	0.80 ⁽²³⁰⁷⁾	0.89 ⁽⁷⁴⁴⁾	0.85 ⁽²¹⁶⁾
VV (Verb, main)	0.74 ⁽¹²⁷⁰⁴⁾	0.82 ⁽⁴⁶⁷²⁾	0.85 ⁽¹⁰⁷³⁾
MM	0.76 ⁽¹¹⁷⁹⁹⁾	0.89 ⁽⁷³³⁾	0.81 ⁽²²³⁾
JC	0.77 ⁽⁷¹²⁾	0.63 ⁽²⁸⁷⁾	0.80 ⁽⁶¹⁾
XSV (Suffix, v.)	0.75 ⁽³⁹⁵⁶⁾	0.85 ⁽¹⁷⁰⁵⁾	0.85 ⁽³⁶⁴⁾
VA (Adjective)	0.73 ⁽⁴⁰²⁸⁾	0.92 ⁽¹⁵⁴⁷⁾	0.81 ⁽³⁹²⁾
NP	0.68 ⁽²²⁶⁰⁾	0.91 ⁽¹⁰¹⁰⁾	0.89 ⁽²⁰¹⁾
NNP	0.65 ⁽³⁶¹⁰⁾	0.47 ⁽³⁴⁷⁶⁾	0.77 ⁽³³⁰⁾
XSA (Suffix, adj.)	0.68 ⁽¹³⁵³⁾	0.71 ⁽³²⁷⁾	0.71 ⁽¹⁴²⁾
VX (Verb, auxiliary)	0.62 ⁽³⁶²⁴⁾	0.64 ⁽¹⁴⁵¹⁾	0.81 ⁽³⁶⁹⁾
XR	0.41 ⁽⁸²⁶⁾	0.67 ⁽³¹⁵⁾	0.49 ⁽⁵²⁾
NR	0.27 ⁽²²⁶⁾	0.78 ⁽⁷³⁾	0.52 ⁽¹⁸⁾

DISCUSSIONS & CONCLUSION

- Overall:** Observation of somewhat reduced performance by the morphological analyzers when applied to L2-Korean data compared to the L1 reference corpus
- By-level:** Identification of asymmetric performance patterns by proficiency level for each analyzer
- By-tag:** Through detailed POS tag analysis, disclosure of low accuracy in essential morphological tags, including predicate (highlighted in yellow) and suffix-related (highlighted in red) tags
- Developing L2 domain-specific model:** Achievement of improved performance in morpheme tokenization and POS tagging for L2-Korean data by integrating L2 data into the training sets
- Future direction:** Emphasis on the potential of carefully designed and validated data-processing pipelines for enhancing computational resources for lesser-studied languages and boosting their performance