

Paying Per-label Attention for Multi-label Extraction from Radiology Reports (Supplementary Material)

Patrick Schrempf^{1,2}, Hannah Watson¹, Shadia Mikhael¹,
Maciej Pajak¹, Matúš Falis¹, Aneta Lisowska¹, Keith W. Muir³,
David Harris-Birtill², and Alison Q. O’Neil^{1,4}

¹ Canon Medical Research Europe, Edinburgh, United Kingdom

² University of St Andrews, United Kingdom

³ Institute of Neuroscience & Psychology, University of Glasgow, United Kingdom

⁴ University of Edinburgh, United Kingdom

patrick.schrempf@eu.medical.canon

Model	#Parameters	Training time [s]	Inference time [s/sample]
BoW + RF	n/a	14 ₁	0.2933 _{0.0040}
Word2Vec	166,524	46 ₈	0.0022 _{0.0001}
CAML [2]	1,021,176	250 ₄₃	0.0099 _{0.0008}
Bi-GRU	2,889,852	111 ₃₂	0.0066 _{0.0003}
Bi-GRU + single attention	3,371,132	120 ₅₅	0.0062 _{0.0003}
Bi-GRU + per-label attention	3,401,852	376 ₁₂₇	0.0109 _{0.0004}
BERT	109,577,596	1115 ₃₂₂	0.0565 _{0.0025}
BioBERT	109,577,596	927 ₁₇₁	0.0575 _{0.0008}
ALARM + softmax	109,458,556	911 ₂₄₃	0.0590 _{0.0013}
ALARM + per-label attention	125,233,276	1448 ₃₇₅	0.0740 _{0.0002}

Table 1: Number of parameters, training time (over 838 samples) and inference time (per sample) for all models. All timings are given as mean_{standard deviation} of 5 runs with different random seeds. The fastest model to train is the random forest model. The Bi-GRU network is significantly faster to train than BERT [1] and ALARM [3] due to the smaller number of parameters. The only model that is faster than the Bi-GRU model is Word2Vec which has a far inferior F1 score. The random forest model is the slowest at inference time because it has n_L models (one model per label) - the inference could be parallelised to improve performance.

Model	Embedding	Data	All	Negative	Uncertain	Positive
Bi-GRU	MIMIC	S	0.584 _{0.022}	0.496 _{0.089}	0.204 _{0.031}	0.642 _{0.012}
Bi-GRU	MIMIC	N-S	0.908 _{0.004}	0.956 _{0.004}	0.427 _{0.058}	0.927 _{0.004}
Bi-GRU	Random	N+S	0.893 _{0.002}	0.962 _{0.008}	0.432 _{0.033}	0.903 _{0.002}
Bi-GRU	MIMIC	N+S	0.921 _{0.003}	0.970 _{0.006}	0.573 _{0.011}	0.932 _{0.004}
ALARM	MIMIC	S	0.569 _{0.028}	0.725 _{0.062}	0.128 _{0.041}	0.531 _{0.028}
ALARM	MIMIC	N-S	0.906 _{0.011}	0.944 _{0.005}	0.532 _{0.087}	0.923 _{0.010}
ALARM	MIMIC	N+S	0.928 _{0.008}	0.965 _{0.004}	0.689 _{0.039}	0.936 _{0.008}

Table 2: Results for our ablation studies showing *micro-averaged* F1 as mean_{standard deviation} of 5 runs with different random seeds (all models are trained with per-label attention). *N* data is the NHS GGC dataset and *S* is the synthetic dataset. “All” combines the classes “negative”, “uncertain” and “positive”. Bold indicates the best model for each metric.

Model	Embedding	Data	All	Negative	Uncertain	Positive
Bi-GRU	MIMIC	S	0.400 _{0.039}	0.504 _{0.050}	0.106 _{0.029}	0.590 _{0.065}
Bi-GRU	MIMIC	N-S	0.551 _{0.024}	0.623 _{0.024}	0.268 _{0.089}	0.761 _{0.026}
Bi-GRU	Random	N+S	0.617 _{0.015}	0.746 _{0.042}	0.360 _{0.054}	0.745 _{0.024}
Bi-GRU	MIMIC	N+S	0.708 _{0.014}	0.796 _{0.027}	0.524 _{0.023}	0.803 _{0.016}
ALARM	MIMIC	S	0.326 _{0.025}	0.607 _{0.039}	0.065 _{0.032}	0.307 _{0.021}
ALARM	MIMIC	N-S	0.534 _{0.041}	0.598 _{0.027}	0.245 _{0.088}	0.758 _{0.038}
ALARM	MIMIC	N+S	0.766 _{0.028}	0.818 _{0.029}	0.661 _{0.061}	0.818 _{0.021}

Table 3: Results for our ablation studies showing *macro-averaged* F1 as mean_{standard deviation} of 5 runs with different random seeds (all models are trained with per-label attention). *N* data is the NHS GGC dataset and *S* is the synthetic dataset. “All” combines the classes “negative”, “uncertain” and “positive”. Bold indicates the best model for each metric.

References

1. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). pp. 4171–4186. Association for Computational Linguistics, Minneapolis, Minnesota (Jun 2019). <https://doi.org/10.18653/v1/N19-1423>
2. Mullenbach, J., Wiegrefe, S., Duke, J., Sun, J., Eisenstein, J.: Explainable prediction of medical codes from clinical text. In: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers). pp. 1101–1111. Association for Computational Linguistics, New Orleans, Louisiana (Jun 2018). <https://doi.org/10.18653/v1/N18-1100>
3. Wood, D., Guilhem, E., Montvila, A., Varsavsky, T., Kiik, M., Siddiqui, J., Kafabadi, S., Gadapa, N., Busaidi, A.A., Townend, M., Patel, K., Barker, G., Ourselin, S., Lynch, J., Cole, J., Booth, T.: Automated Labelling using an Attention model for Radiology reports of MRI scans (ALARM). In: Medical Imaging with Deep Learning (2020), <https://openreview.net/forum?id=UFnWZTbM5t>