

Training Compact Models for Low Resource Entity Tagging using Pre-trained Language Models





Motivation

- Low-resource named entity recognition (NER) is a challenging Natural Language Processing (NLP) task [1] yet widely used for Information Extraction
- In most cases there is a lot of unlabeled data and almost no labeled
- Poor model performance can be fixed by:
- Pre-trained contextual Language Models (LM) were shown to improve many NLP tasks
- Pre-trained LMs are also good for training models with scarcely labeled/low resource data

Pre-trained LMs impose:

Approach

- Knowledge Distillation [2]
 - Teacher-Student training setup
 - Pre-trained LM teacher (BERT)
 - Compact student (~3M params.)
- Pseudo-labeling [3] 2.
 - Utilize abundance of unlabeled data

Tagging more data

- With expert DS/Linguist/Annotator
- Sometimes cannot be fixed because of operational constraints
- Heavy memory
- Heavy compute
- A challenge to deploy in production
- Large model generates *pseudo* entities for the compact model



Experimental Setup

Models

- Teacher BERT [4] (110M/340M params.)
- Compact LSTM-CNN [5] with Softmax/CRF classifier, GloVe embedding (3M params.)

Low-resource Dataset Simulation

- CoNLL 2003 [6] English (PER, ORG, DATE, MISC)
- Randomly sample labeled training sets
- Train set sizes: 150/300/750/1500/3000
- Use remaining training data as unlabeled examples



 $L_{task} = \begin{cases} \mathsf{CrossEntropy}(\hat{y}, y) & labeled example \\ \mathsf{CrossEntropy}(\hat{y}, \hat{y}_{teacher}) & unlabeled example \end{cases}$ Training Steps Fine-tune BERT on labeled data $L_{distillation} = KL(logits_{teacher} || logits_{compact})$ Train compact model using distillation + pseudolabeling loss $Loss = \alpha \cdot L_{task} + \beta \cdot L_{distillation}$ $\alpha + \beta = 1.0$ **Compact Model Accuracy** Inference Speed vs. Teacher Model BERT-base as teacher CPU backend* GPU backend* **BERT-large BERT-base BERT-large BERT-base** 0.85 classifier CRF Softmax CRF Softmax Softmax CRF classifier CRF Batch=1 3.3x 4.3x 8.1x 10.6x Batch=1 0.8x 2.6x 1.5x 0.80 32 33.7x 85.2x 100.4x 28.6x 32 1.3x 5.1x 3.5x ^년 0.75 64 64 45.2x 123.8x 40.0x 109.5x 2.3x 7.2x 6.3x 128 128 49.9x 55.6x 123.6x 137.8x 3.9x 9.7x 10.8x 0.70 **BERT-base** Processor E5-2699A v4 @ 2.40GHz; 251GB RAM; GPU: Titan Xp 12GB; OS: Ubuntu 16.04.1 (4.15.0-50- generic); Tensorflow 1.12 and PyTorch 1.0 Compact+Softmax Compact+CRF DistCompact+Softmax 0.65 DistCompact+CRF Summary





- Compact models have major performance advantage vs. pre-trained LMs in low-resource scenarios
- Teacher models are qualitatively better with more labeled data
- We are exploring:
 - Additional NLP tasks
 - Additional Transformer based pre-trained LM
 - Non-recurrent compact models for fast inference
- Code available in NLP Architect library



peter.izsak@intel.con

Softmax

5.0x

13.1x

19.6x

27.0x

[1] B. Zhang, X. Pan, T. Wang, A. Vaswani, H. Ji, K. Knight, and D. Marcu. Name tagging for low-resource incident languages based on expectation-driven learning. In HLT-NAACL, 2016. [2] G. E. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. ArXiv, abs/1503.02531, 2015.

[4] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. ArXiv, abs/1810.04805, 2018. [5] X. Ma and E. H. Hovy. End-to-end sequence labeling via bi-directional lstm-cnns-crf. ArXiv, abs/1603.01354, 2016.

[3] D.-H. Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep [6] E. F. T. K. Sang and F. D. Meulder. Introduction to the conll-2003 shared task: Languageneural networks. 2013. independent named entity recognition. ArXiv, cs.CL/0306050, 2003.

The 5th Workshop on Energy Efficient Machine Learning and Cognitive Computing @ NeurIPS 2019