

BERTifying the Hidden Markov Model for Multi-Source Weakly Supervised Named Entity Recognition

Yinghao Li¹, Pranav Shetty¹, Lucas Liu¹, Chao Zhang¹, and Le Song²

¹ Georgia Institute of Technology, Atlanta, USA

² Mohamed bin Zayed University of Artificial Intelligence, Abu Dhabi, United Arab Emirates
{yinghaoli, pranav.shetty, lucasliu, chaozhang}@gatech.edu
le.song@mbzuai.ac.ae

Abstract

We study the problem of learning a named entity recognition (NER) tagger using noisy labels from multiple weak supervision sources. Though cheap to obtain, the labels from weak supervision sources are often incomplete, inaccurate, and contradictory, making it difficult to learn an accurate NER model. To address this challenge, we propose a conditional hidden Markov model (CHMM), which can effectively infer true labels from multi-source noisy labels in an unsupervised way. CHMM enhances the classic hidden Markov model with the contextual representation power of pre-trained language models. Specifically, CHMM learns token-wise transition and emission probabilities from the BERT embeddings of the input tokens to infer the latent true labels from noisy observations. We further refine CHMM with an alternate-training approach (CHMM-ALT). It fine-tunes a BERT-NER model with the labels inferred by CHMM, and this BERT-NER’s output is regarded as an additional weak source to train the CHMM in return. Experiments on four NER benchmarks from various domains show that our method outperforms state-of-the-art weakly supervised NER models by wide margins.

1 Introduction

Named entity recognition (NER), which aims to identify named entities from unstructured text, is an information extraction task fundamental to many downstream applications such as event detection (Li et al., 2012), relationship extraction (Bach and Badaskar, 2007), and question answering (Khalid et al., 2008). Existing NER models are typically supervised by a large number of training sequences, each pre-annotated with token-level labels. In practice, however, obtaining such labels could be prohibitively expensive. On the other hand, many domains have various knowledge resources such as

knowledge bases, domain-specific dictionaries, or labeling rules provided by domain experts (Farkiotou et al., 2000; Nadeau and Sekine, 2007). These resources can be used to match a corpus and quickly create large-scale noisy training data for NER from multiple views.

Learning an NER model from multiple weak supervision sources is a challenging problem. While there are works on distantly supervised NER that use only knowledge bases as weak supervision (Mintz et al., 2009; Shang et al., 2018; Cao et al., 2019; Liang et al., 2020), they cannot leverage complementary information from multiple annotation sources. To handle multi-source weak supervision, several recent works (Nguyen et al., 2017; Safranchik et al., 2020; Lison et al., 2020) leverage the hidden Markov model (HMM), by modeling true labels as hidden variables and inferring them from the observed noisy labels through unsupervised learning. Though principled, these models fall short in capturing token semantics and context information, as they either model input tokens as one-hot observations (Nguyen et al., 2017) or do not model them at all (Safranchik et al., 2020; Lison et al., 2020). Moreover, the flexibility of HMM is limited as its transitions and emissions remain constant over time steps, whereas in practice they should depend on the input words.

We propose the conditional hidden Markov model (CHMM) to infer true NER labels from multi-source weak annotations. CHMM conditions the HMM training and inference on BERT by predicting token-wise transition and emission probabilities from the BERT embeddings. These token-wise probabilities are more flexible than HMM’s constant counterpart in modeling how the true labels should evolve according to the input tokens. The context representation ability they inherit from BERT also relieves the Markov constraint and expands HMM’s context-awareness.

Further, we integrate CHMM with a supervised BERT-based NER mode with an alternate-training method (CHMM-ALT). It fine-tunes BERT-NER with the denoised labels generated by CHMM. Taking advantage of the pre-trained knowledge contained in BERT, this process aims to refine the denoised labels by discovering the entity patterns neglected by all of the weak sources. The fine-tuned BERT-NER serves as an additional supervision source, whose output is combined with other weak labels for the next round of CHMM training. CHMM-ALT trains CHMM and BERT-NER alternately until the result is optimized.

Our contributions include:

- A multi-source label aggregator CHMM with token-wise transition and emission probabilities for aggregating multiple sets of NER labels from different weak labeling sources.
- An alternate-training method CHMM-ALT that trains CHMM and BERT-NER in turn utilizing each other’s outputs for multiple loops to optimize the multi-source weakly supervised NER performance.
- A comprehensive evaluation on four NER benchmarks from different domains demonstrates that CHMM-ALT achieves a 4.83 average F1 score improvement over the strongest baseline models.

The code and data used in this work are available at github.com/Yinghao-Li/CHMM-ALT.

2 Related Work

Weakly Supervised NER There have been works that train NER models with different weak supervision approaches. *Distant supervision*, a specific type of weak supervision, generates training labels from knowledge bases (Mintz et al., 2009; Yang et al., 2018; Shang et al., 2018; Cao et al., 2019; Liang et al., 2020). But such a method is limited to one source and falls short of acquiring supplementary annotations from other available resources. Other works adopt multiple additional labeling sources, such as heuristic functions that depend on lexical features, word patterns, or document information (Nadeau and Sekine, 2007; Ratner et al., 2016), and unify their results through multi-source *label denoising*. Several multi-source weakly supervised learning approaches are designed for sentence classification (Ratner et al.,

2017, 2019; Ren et al., 2020; Yu et al., 2020). Although these methods can be adapted for sequence labeling tasks such as NER, they tend to overlook the internal dependency relationship between token-level labels during the inference. Fries et al. (2017) target the NER task, but their method first generates candidate named entity spans and then classifies each span independently. This independence makes it suffer from the same drawback as sentence classification models.

A few works consider label dependency while dealing with multiple supervision sources. Lan et al. (2020) train a BiLSTM-CRF network (Huang et al., 2015) with multiple parallel CRF layers, each for an individual labeling source, and aggregate their transitions with confidence scores predicted by an attention network (Bahdanau et al., 2015; Luong et al., 2015). HMM is a more principled model for multi-source sequential label denoising as the true labels are implicitly inferred through unsupervised learning without deliberately assigning any additional scores. Following this track, Nguyen et al. (2017) and Lison et al. (2020) use a standard HMM with multiple observed variables, each from one labeling source. Safranchik et al. (2020) propose linked HMM, which differs from ordinary HMM by introducing unique linking rules as an adjunct supervision source additional to general token labels. However, these methods fail to utilize the context information embedded in the tokens as effectively as CHMM, and their NER performance is further constrained by the Markov assumption.

Neuralizing the Hidden Markov Model Some works attempt to neuralize HMM in order to relax the Markov assumption while maintaining its generative property (Kim et al., 2018). For example, Dai et al. (2017) and Liu et al. (2018) incorporate recurrent units into the hidden semi-Markov model (HSMM) to segment and label high-dimensional time series; Wiseman et al. (2018) learn discrete template structures for conditional text generation using neuralized HSMM. Wessels and Omlin (2000) and Chiu and Rush (2020) factorize HMM with neural networks to scale it and improve its sequence modeling capacity. The work most related to ours leverages neural HMM for sequence labeling (Tran et al., 2016). CHMM differs from neural HMM in that the tokens are treated as a dependency term in CHMM instead of the observation in neural HMM. Besides, CHMM is trained with generalized EM, whereas neural HMM opti-

companying BERT-NER models are identical to those described in § 5.1. The results in the table suggest that the performance improvement obtained by using alternate-training on the label aggregators is stable and generalizable to any other models yet to be proposed.

6 Conclusion

In this work, we present CHMM-ALT, a multi-source weakly supervised approach that does not depend on manually labeled data to learn an accurate NER tagger. It integrates a label aggregator—CHMM and a supervised model—BERT-NER together into an alternate-training procedure. CHMM conditions HMM on BERT embeddings to achieve greater flexibility and stronger context-awareness. Fine-tuned with CHMM’s prediction, BERT-NER discovers patterns unobserved by the weak sources and complements CHMM. Training these models in turn, CHMM-ALT uses the knowledge encoded in both the weak sources and the pre-trained BERT model to improve the final NER performance. In the future, we will consider imposing more constraints on the transition and emission probabilities, or manipulating them according to sophisticated domain knowledge. This technique could be also extended to other sequence labeling tasks such as semantic role labeling or event extraction.

Acknowledgments

This work was supported by ONR MURI N00014-17-1-2656, NSF III-2008334, Kolon Industries, and research gifts from Google and Amazon. In addition, we would like to thank Yue Yu for his insightful suggestions for this work.

References

Nguyen Bach and Sameer Badaskar. 2007. A review of relation extraction. *Literature review for Language and Statistics II*, 2:1–15.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. [Neural machine translation by jointly learning to align and translate](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. [SciBERT: A pretrained language model for scientific text](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the*

9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3615–3620, Hong Kong, China. Association for Computational Linguistics.

Steven Bird and Edward Loper. 2004. [NLTK: The natural language toolkit](#). In *Proceedings of the ACL Interactive Poster and Demonstration Sessions*, pages 214–217, Barcelona, Spain. Association for Computational Linguistics.

Yixin Cao, Zikun Hu, Tat-seng Chua, Zhiyuan Liu, and Heng Ji. 2019. [Low-resource name tagging learned with weakly labeled data](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 261–270, Hong Kong, China. Association for Computational Linguistics.

Justin Chiu and Alexander Rush. 2020. [Scaling hidden Markov language models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1341–1349, Online. Association for Computational Linguistics.

Hanjun Dai, Bo Dai, Yan-Ming Zhang, Shuang Li, and Le Song. 2017. [Recurrent hidden semi-markov model](#). In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Rezarta Islamaj Dogan, Robert Leaman, and Zhiyong Lu. 2014. [NCBI disease corpus: A resource for disease name recognition and concept normalization](#). *J. Biomed. Informatics*, 47:1–10.

Dimitra Farmakiotou, Vangelis Karkaletsis, John Koutisias, George Sigletos, Constantine D. Spyropoulos, and Panagiotis Stamatopoulos. 2000. Rule-based named entity recognition for greek financial texts. In *Proceedings of the Workshop on Computational Lexicography and Multimedia Dictionaries (COMLEX 2000)*, pages 75–78.

Jason A. Fries, Sen Wu, Alexander Ratner, and Christopher Ré. 2017. [Swellshark: A generative model for biomedical named entity recognition without labeled data](#). *CoRR*, abs/1704.06360.

Zhiheng Huang, Wei Xu, and Kai Yu. 2015. [Bidirectional LSTM-CRF models for sequence tagging](#). *CoRR*, abs/1508.01991.

- Mahboob Alam Khalid, Valentin Jijkoun, and Maarten De Rijke. 2008. The impact of named entity normalization on information retrieval for question answering. In *Proceedings of the IR Research, 30th European Conference on Advances in Information Retrieval, ECIR'08*, pages 705–710, Berlin, Heidelberg. Springer-Verlag.
- Yoon Kim, Sam Wiseman, and Alexander M. Rush. 2018. A tutorial on deep latent variable models of natural language. *CoRR*, abs/1812.06834.
- Ouyu Lan, Xiao Huang, Bill Yuchen Lin, He Jiang, Liyuan Liu, and Xiang Ren. 2020. Learning to contextually aggregate multi-source supervision for sequence labeling. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2134–2146, Online. Association for Computational Linguistics.
- Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2019. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.
- Chenliang Li, Aixin Sun, and Anwitaman Datta. 2012. Twevent: Segment-based event detection from tweets. In *Proceedings of the 21st ACM International Conference on Information and Knowledge Management, CIKM '12*, pages 155–164, New York, NY, USA. Association for Computing Machinery.
- Jiao Li, Yueping Sun, Robin J. Johnson, Daniela Sciaky, Chih-Hsuan Wei, Robert Leaman, Allan Peter Davis, Carolyn J. Mattingly, Thomas C. Wieggers, and Zhiyong Lu. 2016. Biocreative V CDR task corpus: a resource for chemical disease relation extraction. *Database J. Biol. Databases Curation*, 2016.
- Chen Liang, Yue Yu, Haoming Jiang, Siawpeng Er, Ruijia Wang, Tuo Zhao, and Chao Zhang. 2020. Bond: Bert-assisted open-domain named entity recognition with distant supervision. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '20*, pages 1054–1064, New York, NY, USA. Association for Computing Machinery.
- Pierre Lison, Jeremy Barnes, Aliaksandr Hubin, and Samia Touileb. 2020. Named entity recognition without labelled data: A weak supervision approach. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1518–1533, Online. Association for Computational Linguistics.
- Hao Liu, Lirong He, Haoli Bai, Bo Dai, Kun Bai, and Zenglin Xu. 2018. Structured inference for recurrent hidden semi-markov model. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI'18*, page 2447–2453. AAAI Press.
- Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Effective approaches to attention-based neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1412–1421, Lisbon, Portugal. Association for Computational Linguistics.
- Mike Mintz, Steven Bills, Rion Snow, and Daniel Jurafsky. 2009. Distant supervision for relation extraction without labeled data. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 1003–1011, Suntec, Singapore. Association for Computational Linguistics.
- David Nadeau and Satoshi Sekine. 2007. A survey of named entity recognition and classification. *Linguisticae Investigationes*, 30(1):3–26. Publisher: John Benjamins Publishing Company.
- An Thanh Nguyen, Byron Wallace, Junyi Jessy Li, Ani Nenkova, and Matthew Lease. 2017. Aggregating and predicting sequence labels from crowd annotations. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 299–309, Vancouver, Canada. Association for Computational Linguistics.
- Maria Pontiki, Dimitris Galanis, John Pavlopoulos, Harris Papageorgiou, Ion Androutsopoulos, and Suresh Manandhar. 2014. SemEval-2014 task 4: Aspect based sentiment analysis. In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 27–35, Dublin, Ireland. Association for Computational Linguistics.
- Lance Ramshaw and Mitch Marcus. 1995. Text chunking using transformation-based learning. In *Third Workshop on Very Large Corpora*.
- Alexander Ratner, Stephen H. Bach, Henry Ehrenberg, Jason Fries, Sen Wu, and Christopher Ré. 2017. Snorkel: Rapid training data creation with weak supervision. *Proc. VLDB Endow.*, 11(3):269–282.
- Alexander Ratner, Braden Hancock, Jared Dunnmon, Frederic Sala, Shreyash Pandey, and Christopher Ré. 2019. Training complex models with multi-task weak supervision. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33:4763–4771.
- Alexander Ratner, Christopher De Sa, Sen Wu, Daniel Selsam, and Christopher Ré. 2016. Data programming: Creating large training sets, quickly. In *Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS'16*, pages 3574–3582, Red Hook, NY, USA. Curran Associates Inc.
- Wendi Ren, Yinghao Li, Hanting Su, David Kartchner, Cassie Mitchell, and Chao Zhang. 2020. Denoising multi-source weak supervision for neural text classification. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3739–3754, Online. Association for Computational Linguistics.

- Esteban Safranchik, Shiyong Luo, and Stephen H. Bach. 2020. [Weakly supervised sequence tagging from noisy rules](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 5570–5578. AAAI Press.
- Jingbo Shang, Liyuan Liu, Xiaotao Gu, Xiang Ren, Teng Ren, and Jiawei Han. 2018. [Learning named entity tagger using domain-specific dictionary](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2054–2064, Brussels, Belgium. Association for Computational Linguistics.
- Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. 2019. [How to fine-tune bert for text classification?](#) *Chinese Computational Linguistics*, page 194–206.
- Christian Thiel. 2008. [Classification on soft labels is robust against label noise](#). In *Proceedings of the 12th International Conference on Knowledge-Based Intelligent Information and Engineering Systems, Part I, KES '08*, pages 65–73, Berlin, Heidelberg. Springer-Verlag.
- Erik F. Tjong Kim Sang and Fien De Meulder. 2003. [Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition](#). In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, pages 142–147.
- Ke M. Tran, Yonatan Bisk, Ashish Vaswani, Daniel Marcu, and Kevin Knight. 2016. [Unsupervised neural hidden Markov models](#). In *Proceedings of the Workshop on Structured Prediction for NLP*, pages 63–71, Austin, TX. Association for Computational Linguistics.
- A. Viterbi. 1967. [Error bounds for convolutional codes and an asymptotically optimum decoding algorithm](#). *IEEE Trans. Inf. Theor.*, 13(2):260–269.
- T. Wessels and Christian W. Omlin. 2000. [Refining hidden markov models with recurrent neural networks](#). In *Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks, IJCNN 2000, Neural Computing: New Challenges and Perspectives for the New Millennium, Como, Italy, July 24-27, 2000, Volume 2*, pages 271–278. IEEE Computer Society.
- Sam Wiseman, Stuart Shieber, and Alexander Rush. 2018. [Learning neural templates for text generation](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3174–3187, Brussels, Belgium. Association for Computational Linguistics.
- Yaosheng Yang, Wenliang Chen, Zhenghua Li, Zhengqiu He, and Min Zhang. 2018. [Distantly supervised NER with partial annotation learning and reinforcement learning](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2159–2169, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Yue Yu, Simiao Zuo, Haoming Jiang, Wendi Ren, Tuo Zhao, and Chao Zhang. 2020. [Fine-tuning pre-trained language model with weak supervision: A contrastive-regularized self-training approach](#).

A Technical Details

A.1 CHMM Training

Following the discussion in § 4.2, we use the *forward-backward* algorithm to calculate the smoothed marginal $\gamma_i^{(t)} \triangleq p(z^{(t)} = i | \mathbf{x}^{(1:T)})$, $i \in \{1, 2, \dots, |\mathcal{L}|\}$, $t \in \{1, 2, \dots, T\}$ and the expected number of transitions $\xi_{i,j}^{(t)} \triangleq p(z^{(t-1)} = i, z^{(t)} = j | \mathbf{x}^{(1:T)})$, $i, j \in \{1, 2, \dots, |\mathcal{L}|\}$.⁸ $|\mathcal{L}|$ is the number of BIO formatted entity labels, which are regarded as hidden states; T is the total number of hidden steps in a sequence, which equals the number of tokens.

Defining $\alpha_i^{(t)} \triangleq p(z^{(t)} = i | \mathbf{x}^{(1:t)})$ and $\beta_i^{(t)} \triangleq p(\mathbf{x}^{(t+1:T)} | z^{(t)} = i)$, $\gamma_i^{(t)}$ and $\xi_{i,j}^{(t)}$ can be represented by α and β using the Bayes' rule and Markov assumption:

$$\begin{aligned} \gamma_i^{(t)} &\triangleq p(z^{(t)} = i | \mathbf{x}^{(1:T)}) \\ &= \frac{p(\mathbf{x}^{(t+1:T)}, z^{(t)} = i | \mathbf{x}^{(1:t)})}{p(\mathbf{x}^{(t+1:T)} | \mathbf{x}^{(1:T)})} \\ &\propto p(z^{(t)} = i | \mathbf{x}^{(1:t)}) p(\mathbf{x}^{(t+1:T)} | z^{(t)} = i) \\ &= \alpha_i^{(t)} \beta_i^{(t)}, \end{aligned} \quad (12)$$

$$\begin{aligned} \xi_{i,j}^{(t)} &\triangleq p(z^{(t-1)} = i, z^{(t)} = j | \mathbf{x}^{(1:T)}) \\ &\propto p(z^{(t-1)} = i | \mathbf{x}^{(1:t-1)}) \\ &\quad p(z^{(t)} = j | z^{(t-1)} = i, \mathbf{x}^{(t:T)}) \\ &\propto p(z^{(t-1)} = i | \mathbf{x}^{(1:t-1)}) p(\mathbf{x}^{(t)} | z^{(t)} = j) \\ &\quad p(\mathbf{x}^{(t+1:T)} | z^{(t)} = j) p(z^{(t)} = j | z^{(t-1)} = i) \\ &= \alpha_i^{(t-1)} \varphi_j^{(t)} \beta_j^{(t)} \Psi_{i,j}^{(t)}. \end{aligned} \quad (13)$$

$\varphi_i^{(t)} \in \mathbb{R}^{|\mathcal{L}|} \triangleq p(\mathbf{x}^{(t)} | z^{(t)} = i)$ is the likelihood of the observation when the hidden state is i (§ 4.2).

Written in the matrix form, (12) and (13) become:

$$\boldsymbol{\gamma}^{(t)} \propto \boldsymbol{\alpha}^{(t)} \odot \boldsymbol{\beta}^{(t)}, \quad (14)$$

$$\boldsymbol{\xi}^{(t)} \propto \boldsymbol{\Psi}^{(t)} \odot (\boldsymbol{\alpha}^{(t-1)} (\boldsymbol{\varphi}^{(t)} \odot \boldsymbol{\beta}^{(t)})^\top), \quad (15)$$

where \odot is the element-wise product. Note that the elements in both $\boldsymbol{\gamma}^{(t)}$ and $\boldsymbol{\xi}^{(t)}$ should sum up to 1.

⁸Same as § 4.2, we omit the dependency term $e^{(1:T)}$.

The Forward Pass The filtered marginal $\alpha_i^{(t)}$ can be computed iteratively:

$$\begin{aligned} \alpha_i^{(t)} &\triangleq p(z^{(t)} = i | \mathbf{x}^{(1:t)}) \\ &= p(z^{(t)} = i | \mathbf{x}^{(t)}, \mathbf{x}^{(1:t-1)}) \\ &\propto p(\mathbf{x}^{(t)} | z^{(t)} = i) p(z^{(t)} = i | \mathbf{x}^{(1:t-1)}) \quad (16) \\ &= \sum_j \varphi_i^{(t)} \Psi_{j,i}^{(t)} \alpha_j^{(t-1)}. \end{aligned}$$

Written in the matrix form, (16) becomes

$$\boldsymbol{\alpha}^{(t)} \propto \boldsymbol{\varphi}^{(t)} \odot (\boldsymbol{\Psi}^{(t)})^\top \boldsymbol{\alpha}^{(t-1)}. \quad (17)$$

We initialize $\boldsymbol{\alpha}$ with $\boldsymbol{\alpha}^{(0)} = \boldsymbol{\pi}$ (§ 4.2) since we have no observation at time step 0. As $\boldsymbol{\alpha}^{(t)}$ is a probability distribution, the elements in it sum up to 1. The calculation of $\boldsymbol{\alpha}$ is the *forward pass*.

The Backward Pass In the same way, we do the *backward pass* to compute the conditional future evidence $\beta_i^{(t)} \triangleq p(\mathbf{x}^{(t+1:T)} | z^{(t)} = i)$:

$$\begin{aligned} \beta_i^{(t-1)} &\triangleq p(\mathbf{x}^{(t+1:T)} | z^{(t)} = j) \\ &= \sum_j p(z^{(t)} = j, \mathbf{x}^{(t)}, \mathbf{x}^{(t+1:T)} | z^{(t-1)} = i) \\ &= \sum_j [p(\mathbf{x}^{(t+1:T)} | z^{(t)} = j) \\ &\quad p(\mathbf{x}^{(t)}, z^{(t)} = j | z^{(t-1)} = i)] \\ &= \sum_j \beta_j^{(t)} \varphi_j^{(t)} \Psi_{i,j}^{(t)}. \end{aligned} \quad (18)$$

In the matrix form, (18) becomes:

$$\boldsymbol{\beta}^{(t-1)} = \boldsymbol{\Psi}^{(t)} (\boldsymbol{\varphi}^{(t)} \odot \boldsymbol{\beta}^{(t)}), \quad (19)$$

whose base case is

$$\begin{aligned} \beta_i^{(T)} &= p(\mathbf{x}^{(T+1:T)} | z^{(T)} = i) = 1, \\ &\quad \forall i \in \{1, \dots, |\mathcal{L}|\}. \end{aligned}$$

A.2 The Maximization step for Unsupervised HMM

For traditional unsupervised HMM, the expected complete data log likelihood is maximized by updating the matrices with the approximated pseudo-statistics. different from CHMM, HMM has constant transition and emission for all time steps, *i.e.*:

$$\boldsymbol{\Psi}^{(1)} = \boldsymbol{\Psi}^{(t)}; \boldsymbol{\Phi}^{(1)} = \boldsymbol{\Phi}^{(t)}; \quad \forall t \in \{2, \dots, T\}.$$

For simplicity, we remove the term t for the transition and emission matrices. Suppose we are updating HMM based on one instance with t starting from 1:

$$\pi_i = \gamma_i^{(1)}; \quad (20)$$

$$\Psi_{i,j} = \frac{\sum_{t=2}^T \xi_{i,j}^{(t)}}{\sum_{t=2}^T \sum_{\ell=1}^{|\mathcal{L}|} \xi_{i,\ell}^{(t)}}; \quad (21)$$

$$\Phi_{i,j,k} = \frac{\sum_{t=1}^T \gamma_i^{(t)} x_{j,k}^{(t)}}{\sum_{t=1}^T \gamma_i^t}. \quad (22)$$

Note that the observation has property $0 \leq x_{j,k}^{(t)} \leq 1$ and $\sum_{j=1}^{|\mathcal{L}|} x_{j,k}^{(t)} = 1$, where $k \in \{1, \dots, K\}$ is the index of the weak labeling source.

B Labeling Source Performance

The weak labeling sources of the CoNLL 2003 dataset come from [Lison et al. \(2020\)](#), whereas [Safranchik et al. \(2020\)](#) provide the sources for the LaptopReview, NCBI-Disease and BC5CDR dataset. For [Safranchik et al. \(2020\)](#)'s labeling sources, we apply a majority voting using their tagging results to the spans detected by their *linking rules* to convert the linking results to token annotations. In consideration of the training time and resource consumption, we only adopt a subset of the labeling sources provided by the authors. The performance of the labeling sources is presented in the tables below.

source name	precision	recall	f1
CoreDictionaryUncased	81.03	41.41	5.48
CoreDictionaryExact	80.69	17.18	28.32
CancerLike	34.88	1.58	3.02
BodyTerms	68.52	3.90	7.38
ExtractedPhrase	97.12	32.03	48.18

Table 4: The performance of the labeling sources used in the NCBI-Disease dataset.

source name	precision	recall	f1
DictCore-Chemical	91.81	29.55	44.7
DictCore-Chemical-Exact	85.88	3.16	6.1
DictCore-Disease	81.57	26.32	39.8
DictCore-Disease-Exact	81.4	1.09	2.16
Organic Chemical	92.67	30.07	45.4
Disease or Syndrome	77.36	11.67	20.28
PostHyphen	84.47	08.07	14.74
ExtractedPhrase	86.8	17.96	29.76

Table 5: The performance of the labeling sources used in the BC5CDR dataset.

source name	precision	recall	f1
CoreDictionary	72.63	51.61	60.34
iStuff	26.67	0.61	1.2
ExtractedPhrase	97.45	29.25	45.0
ConsecutiveCapitals	35.29	0.92	1.8

Table 6: The performance of the labeling sources used in the LaptopReview dataset.

source name	precision	recall	f1
BTC+c	61.56	46.35	52.88
SEC+c	39.54	24.59	30.32
core_web_md+c	69.53	60.04	64.44
crunchbase_cased	38.26	5.59	9.76
crunchbase_uncased	37.88	6.2	10.66
doc_majority_cased	65.81	40.21	49.92
doc_majority_uncased	61.69	40.17	48.66
full_name_detector	87.79	11.33	20.06
geo_cased	68.16	15.35	25.06
geo_uncased	65.1	18.89	29.28
misc_detector	85.14	21.51	34.34
wiki_cased	75.27	32.65	45.54
wiki_uncased	72.26	35.61	47.7

Table 7: The performance of the labeling sources used in the CoNLL 2003 dataset.

Please refer to [Lison et al. \(2020\)](#) for the information about the construction of the labeling sources on the CoNLL 2003 dataset; please refer to [Safranchik et al. \(2020\)](#) for the labeling sources on other three datasets.

C Hyper-Parameters

The experiments are conducted on one GeForce RTX 2080 Ti GPU. For NCBI-Disease, BC5CDR and LaptopReview datasets, CHMM is pre-trained for 5 epochs and trained for 20 epochs. The learning rates for these three datasets are 5×10^{-4} , 10^{-3} and 10^{-4} , respectively, and the batch sizes are 64, 64 and 128. In phase I, BERT-NER is trained with the default learning rate (5×10^{-5}) for 100 epochs. The batch sizes are 8, 8, and 48, respectively. Note that for LaptopReview, the maximum length limitation of BERT-NER is set to 128 whereas the limitation is 512 for the other two datasets. In phase II, we use half the learning rate with 20 epochs for each loop.

For CoNLL 2003, CHMM has the same number of training epochs as for other datasets. The batch size is 32, and the learning rate is 10^{-5} . BERT-NER has a maximum sequence length of 256. It is trained for 15 epochs in phase I and 5 epochs in phase II. Other hyper-parameters are identical to other BERT-NER models'.