# Weakly Supervised Named Entity Recognition for Historical Texts

**Marco Sorbi[1,2], Laurent Moccozet[2], Stephane Marchand-Maillet[2]**

[1]Research Institute for Statistics and Information Science, University of Geneva, Switzerland
[2]Centre Universitaire d'Informatique, University of Geneva, Switzerland
**Correspondence:** Marco.Sorbi@unige.ch

## Abstract

Named Entity Recognition has emerged as a critical task in natural language processing, particularly for extracting meaningful information from unstructured text. Although traditional approaches rely heavily on large annotated datasets, recent advances have explored weak supervision techniques to address the limitations of resource-intensive annotation processes. Historical texts provide unique challenges to this task because of their linguistic peculiarities, and several approaches exist to address texts of this domain in a supervised way, but they involve lengthy manual annotations of the documents of interest by domain experts. To address this issue, this paper explores how recent weakly supervised NER techniques can be adapted to historical texts, analyzing their suitability for this domain. The experiments show that domain-specific architectures can be effectively trained on low-resource corpora with weak supervision over a small set of entity labels. Using only 10% of the annotations, the performance of these architectures remains above 80% of the supervised quality in terms of F1-Score.

## 1 Introduction

Named Entity Recognition (NER) is a fundamental task in Natural Language Processing to identify and classify named entities from unstructured text into predefined categories such as people, organizations, and locations, and serves as an important step in the structuring of textual data (Lauriola et al., 2021). General purpose NER methods range from traditional rule-based approaches to modern deep learning architectures, mainly transformer-based (Keraghel et al., 2024). Various domains exist where ad hoc methods have been studied to address domain-specific peculiarities, such as biomedicine (Lauriola et al., 2021) and history (Ehrmann et al., 2023).

Historical documents present distinct challenges compared to modern texts, ranging from an archaic language with non-standardized orthography (Santini et al., 2025), to degraded scripts and complex document layouts (Ehrmann et al., 2020). The lack of large annotated corpora in this low-resource context makes it difficult to design and train specific NER architectures (Novotny et al., 2023) that are required for effective NER in historical documents, with a significant impact on digital humanities and cultural heritage activities.

When manually annotated data are scarce, weakly supervised NER (WS-NER) has emerged to address the lack of training data. One common WS-NER framework is distantly-supervised NER (DS-NER), which takes advantage of knowledge bases or dictionaries to automatically generate labels to use as supervision signals (Fang et al., 2021). However, DS-NER faces significant challenges due to noisy and incomplete labels caused by limited dictionary coverage and imperfect distant annotation (Zhou et al., 2022). Previous works have studied multiple approaches to DS-NER, which we will present in Section 2, which have the potential to progress NER in the historical domain given its noise context and lack of extensive labeled datasets.

This paper aims to study WS-NER in the historical domain, exploiting existing DS-NER techniques for modern texts, to evaluate their effectiveness in the domain, and use them to adapt architectures from supervised Historical NER studies to the WS-NER task[1]. This work is motivated by the feasibility of creating a partial domain-specific knowledge base relative to labeling a complete domain-specific document to train a supervised NER system. For this purpose, we identify two research questions to address:

- **RQ1**: Varying the size of the dictionaries, how do DS-NER techniques that are designed for modern texts perform on historical documents with domain-specific dictionaries? How do

---

[1] Code is available at www.github.com/msorbi/hwsner

they compare to domain-specific supervised NER methods in WS-NER and fully supervised NER settings?

- **RQ2**: Can we exploit these DS-NER techniques to adapt domain-specific supervised NER architectures to WS-NER?

The remainder of this paper is structured as follows. In Section 2, we discuss and organize a selection of related works concerning DS-NER and supervised NER in historical documents. Then, Section 3 presents the datasets used in this work to evaluate the models. Section 4 describes the experiments carried out and presents an overview of the metrics used for evaluation in this work. The results are presented in Section 5 and discussed in Section 6, where the research questions are also answered. Finally, Section 7 presents conclusions and suggestions for future work.

## 2 Related Work

The positioning of this paper is at the intersection of weak supervision and the historical domain of Named Entity Recognition. The presentation of an overview of related work on the two topics is therefore appropriate.

### 2.1 Distantly-supervised Named Entity Recognition

DS-NER is a WS-NER framework that relies on the availability of dictionaries or knowledge bases instead of labeled sequences as training supervision signal. It avoids the need for large-scale human annotation, but introduces noise due to incomplete and inaccurate annotations (Zhang et al., 2021). Early research sees the development of AutoNER (Shang et al., 2018), which introduced a Tie-or-Break tagging scheme with a Fuzzy-LSTM-CRF architecture designed to be robust to noisy distant supervision. We divide subsequent work into two main branches: Positive-Unlabeled Learning (PUL) and Self-Supervised Learning (SSL).

**Positive-Unlabeled Learning** It formulates DS-NER as training on positively labeled data, created using the knowledge base, and unlabeled data. Contributions using this approach include:

- AdaPU (Peng et al., 2019), which firstly introduced PUL in DS-NER, designing a training algorithm that can unbiasedly estimate the task loss as if there were fully labeled data.

- Conf-MPU (Zhou et al., 2022), which extended PUL in DS-NER to a multi-class setting and introduced confidence-based risk estimation.

- CuPUL (Li et al., 2025), which adds curriculum learning to PUL to stabilize training and weaken the impact of noise.

**Self-Supervised Learning** Used as a training stage, it allows reducing the impact of noisy labels to works like:

- BOND (Liang et al., 2020), which first adapts Pre-trained Language Models to noisy distant labels, and then applies teacher-student self-training for refinement.

- SCDL (Zhang et al., 2021), which trains two teacher-student networks to jointly denoise labels.

- CENSOR (Si et al., 2024), which introduces uncertainty-aware teacher learning, to reduce reliance on miscalibrated high-confidence labels, and student-student label sharing, to mitigate error propagation.

**Other approaches** There are other recent works using different approaches, for example:

- SANTA (Si et al., 2023), addressing inaccurate and incomplete annotations separately, with a memory-smoothed focal loss and a noise-tolerant loss.

- MProto (Wu et al., 2023), employing a prototype network to capture intra-class variance, and formulating token-prototype assignment as an optimal transport problem.

### 2.2 Historical Named Entity Recognition

In the domain of historical documents, NER becomes substantially more difficult, as it faces multiple challenges due to the linguistic characteristics of these texts, which the survey by Ehrmann et al. (2023) identifies in language dynamics, noisy input, and lack of resources. Analyzing methodologies from rule-based to deep learning, the survey emphasizes the need for models robust to multilingualism, spelling variation, and sparse data. When dealing with digitized documents, non-standard layout and Optical Character Recognition (OCR) errors add a layer of noise that can significantly degrade performance (Kettunen et al., 2017).

We can distinguish the application of different techniques for NER in this domain.

**Rule-based**  Early research mainly employs rule-based systems tailored to documents of interest, using lexical heuristics, gazetteers, and token-level rules, in order to extract people and places from British parliamentary records (Grover et al., 2008), for example. This approach is transparent, but typically struggles to face the challenge of integrating variability in the rules when increasing the corpus size.

**Machine learning**  It enables the use of more flexible feature-based techniques that commonly rely on Hidden Markov Models, Support Vector Machines, Conditional Random Fields (CRF), and decision trees (Yadav and Bethard, 2018). In historical context, Torres Aguilar et al. (2016) applied a CRF on latin charters, obtaining good results based on a large number of features. These methods reduce the effect of variability, but they require a large number of hand designed features.

**Deep learning**  These systems, often relying on sequence labeling and transformer-based models such as Long Short-Term Memory (LSTM) and Bidirectional Encoder Representations from Transformers (BERT) for automated feature learning, and CRF for classification, have further reduced the effect of variability. Among deep learning works in the historical domain,

- Boros et al. (2020) propose a hierarchical transformer stack for historical German and French datasets, specifically to mitigate OCR noise in digitized historical texts,

- Blouin et al. (2021) investigate transfer learning from modern to historical domains, studying annotation effort, domain mismatch, and pre-training data selection for historical NER,

- Torres Aguilar (2022) creates a human-labeled dataset of multilingual medieval charters, and addresses multilingual NER by combining stacked embeddings with BERT-based models fine-tuned on the dataset,

- Novotny et al. (2023) create a small dataset with sentences containing known entities from late medieval European texts, and propose a bootstrapping annotation pipeline to build larger copora of the texts.

However, these methods rely on heavily annotated datasets to learn features, which are rarely available.

**Few-shot learning**  It exploits Large Lanugage Models (LLMs) through prompting in order to label the historical texts without relying on annotated corpora. Among these works in the historical domain,

- Hiltmann et al. (2025) prompt LLMs with historical context, outperforming established frameworks on historical texts,

- Zhang and Colavizza (2025) use few-shot learning of LLMs on the HIPE dataset (Ehrmann et al., 2020).

This work focuses on deep learning systems to reduce the amount of annotations required to train them. Specifically, we study the applicability of weak supervision techniques to train deep learning models on the historical NER, studying the impact of the amount of annotated data on the quality of results. Even if we can consider the approach of Novotny et al. (2023) to belong to the WS-NER paradigm, to the knowledge of the authors, there is no extensive study of WS-NER or DS-NER in the historical domain.

## 3  Dataset

Our experiments are based on the dataset introduced by Torres Aguilar (2022), consisting of a human-labeled medieval NER dataset whose ancient French charters come from:

- Corpus de la Bourgogne du Moyen Âge ("Corpus of Burgundy in the Middle Ages", *CBMA*) - Cartulary of the city of Arbois (Magnani, 2020): a municipal cartulary commissioned in 1384, containing public issues such as military services and war costs, taxes and customs, or lawsuits in court.

- Diplomata Belgica ("Belgian Diplomatic Sources", *CDBE* – de Hemptinne, Thérèse and Deploige, Jeroen and Kupper, Jean-Louis and Prevenier, Walter, 2015): a database published by the Belgian Royal Historical Commission, containing French charters dated the 13th century, containing legal actions concerning individuals, corporations, and private affairs.

- HOME History of Medieval Europe (*HOME*) - Alcar (Stutzmann et al., 2021): it contains cartularies dated between the 12th and 14th centuries, reporting donations, exchanges, and other legal acts.

|            | CBMA  | CDBE   | HOME   |
|------------|-------|--------|--------|
| Train      | 38658 | 235643 | 114640 |
| Validation | 896   | 18510  | 18510  |
| Test       | 8133  | 56081  | 18554  |

Table 1: Number of tokens per each dataset split.

|      | CBMA | CDBE | HOME |
|------|------|------|------|
| PERS | 652  | 4118 | 925  |
| LOC  | 347  | 2662 | 922  |

Table 2: Number of distinct training entities for each type and source, before sampling.

The dataset contains annotations for entities in the Person (PERS) and Location (LOC) classes, and the sizes for each source are reported in Table 1.

Separately for each source, to adapt the dataset to a WS-NER setting, we extracted all the tagged entities from the train split, and randomly sampled the entities to make the dictionaries for each type. Then, we used the sampled dictionaries to create the noisy labels for the train splits of the mentioned sources via string matching. Table 2 shows the dimensions of each dictionary. The validation and test splits remained unchanged.

## 4 Experimental Setup

To address RQ1, we tested DS-NER techniques belonging to the Positive-Unlabeled and Self-Supervised Learning frameworks, plus the two specific works mentioned in Section 2 and the Stacked Embeddings supervised architecture (Torres Aguilar, 2022), and compared their performances on historical French texts in two settings:

- Fully supervised NER, using the original human-labeled train splits, to analyze the best supervision condition, and

- WS-NER, varying the sampling size for the dictionaries according to Section 3, to see its impact on the performances. We used sampling sizes from 20% to 100%, increasing by 20%, that are commonly used in previous works (Zhou et al., 2022; Wu et al., 2023), and 10% to test with increased noise conditions,

training every method on the train split of each setting and sample size independently. To better summarize the results, we will refer with "small

dictionaries" to the 10% and 20% samples, with "medium dictionaries" to the 40% and 60% samples, and with "large dictionaries" to the 80% and 100% samples.

Then, to address RQ2, we used the best performing DS-NER technique to train

1. XLM-RoBERTa base (Conneau et al., 2020), a large multilingual masked language model, and

2. Stacked Embeddings (Torres Aguilar, 2022), a domain-specific architecture for supervised NER in medieval texts,

on the WS-NER task, using the same settings as in the first part, and checked whether DS-NER can generalize their performance from the supervised task.

**Evaluation**  We now discuss the evaluation metrics used to assess NER performance in historical documents. According to the works presented in Section 2, the typical evaluation metric for NER is the F1-score, usually presented alongside precision and recall to show their impact. Novotny et al. (2023) use the $F_{0.25}$-score, which significantly gives more weight to precision than to recall, but otherwise the two metrics are averaged in a balanced way, that is, using the F1-score.

In WS-NER, training supervision is based on incomplete labels, which are characterized by high precision but low recall (Zhou et al., 2022). This effect is amplified in historical documents by orthographic variations that complicate entity identification. As DS-NER training algorithms aim to improve the recall at the expense of precision, we continue to use the common F1-score as a balanced measure, giving importance to the recall metric, too. We compute these metrics using seqeval (Nakayama, 2018) default token-level mode, in order to grant partial scores to entities that are correctly identified and classified, even if their boundary is not correct, and mitigate the problem of annotation granularity (Ehrmann et al., 2020).

If multiple classes are present, as in our case with PERS and LOC, micro averaging is commonly used for global measures (Shang et al., 2018; Torres Aguilar, 2022) to ensure that frequent classes have an appropriate influence on the final results.

## 5 Results

We analyze here the test performances, in terms of micro-averaged F1-Score, Precision, and Recall,

obtained by the various methods under fully supervised NER and different weak supervision settings. For each technique, we micro-average measures among classes and then, for each setting, we average among datasets and report them in percentage in Figure 1. Detailed measures are reported in Appendix A. For comparison, Stacked Embeddings (StEmb), a domain-specific architecture for supervised NER in medieval texts (Torres Aguilar, 2022), is also reported along with dictionary matching and the CuPUL DS-NER technique paired with XLM-RoBERTa (XLM) and Stacked Embeddings as underlying classifier. We can notice that:

- Despite having a precision among the best in every setting, if we consider recall and F1-score, the Stacked Embeddings architecture is the best performer with full supervision and large dictionaries, but its performances drop quickly when the sample size is decreased.

- If trained under full supervision, the best DS-NER techniques are BOND (Liang et al., 2020) and SANTA (Si et al., 2023), with an F1-score close to the baseline supervised Stacked Embeddings architecture.

- SANTA maintains good recall with large dictionaries, but if we consider the F1-score, BOND remains the best DS-NER technique with large and medium dictionaries, overcoming the supervised baseline as well for medium dictionaries, and is the best DS-NER technique in every setting in terms of precision.

- With small dictionaries, the best technique in terms of F1-score is CuPUL (Li et al., 2025), which has the best recall even with medium dictionaries.

- Considering all three metrics, CuPUL is the least affected by the noise introduced with distant labels, with a performance that only slightly changes on the different settings. Conf-MPU (Zhou et al., 2022), which is the other PUL-based technique, and MProto (Wu et al., 2023) also have stable but lower performances.

- Except for the simpler BOND, the SSL techniques obtained generally unstable performances across settings, especially in terms of recall and F1-score.

Because CuPUL is the best performing technique in low supervision settings, and the least affected by the noise of distant labels, we use it to train (1) XLM-RoBERTa (Conneau et al., 2020) and (2) Stacked Embeddings (Torres Aguilar, 2022) on the WS-NER task. This is done by replacing the original RoBERTa classifier in CuPUL. Despite that they perform slightly worse than the baseline with full supervision and large dictionaries, we can see once again in Figure 1 that CuPUL is little affected by the noise of distant labels. Although XLM-RoBERTa (1) still has some noise-related problems, especially with small dictionaries, CuPUL allows the ad-hoc classifier (2) to have performances in the WS-NER settings that are significantly closer to that of the fully supervised task for which it has been designed.

## 6 Discussion

The results of our experiments show that weak supervision can achieve F1-Scores within 20% of supervised Named Entity Recognition performance on historical texts using only 10% of the annotations, making them particularly appealing for this context where large annotated datasets are often unavailable. To mitigate the linguistic challenges that characterize NER in historical texts, the experiments demonstrate that DS-NER techniques can effectively adapt domain-specific NER models to low-resource historical corpora, to achieve performance levels that are close to fully supervised methods even with additional noise introduced with distant labels. The amount of noise in the distant labels, controlled via the dictionary size, has a significant impact on the final performance, causing noticeable drops in most of the techniques. Nevertheless, CuPUL and MProto have shown good robustness, maintaining consistent performance across different dictionary sizes, and revealing particularly suitable for this scenario. Considering these results, we can now answer the research questions formulated in Section 1.

### 6.1 RQ1: Validity of modern text DS-NER techniques for ancient texts

The results reported in Figure 1 show that most DS-NER techniques suffer from performance drops when they are applied to NER in ancient texts, especially in low supervision settings with small dictionaries. CuPUL (Li et al., 2025), BOND (Liang et al., 2020), and SANTA (Si et al., 2023), belong-

Figure 1: Test results of Fully supervised NER and Weakly supervised NER varying the dictionary sample size. The Stacked Embeddings (StEmb) supervised architecture, dictionary matching, and the CuPUL DS-NER technique paired with XLM-RoBERTa (XLM) and Stacked Embeddings as underlying classifier are reported using solid lines.

ing to three different categories as identified in Section 2, are able to achieve good performance with large dictionaries, but still suffer substantial drops when switching to smaller dictionaries, especially the latter two techniques. CuPUL is the most robust to the noise of low supervision, thanks to a high recall as opposed to the other techniques.

## 6.2 RQ2: Adaptation of domain-specific supervised NER architectures using modern text DS-NER

Based on the answer to RQ1, we selected CuPUL to carry out further experiments on the adaptability to the WS-NER task in historical texts. Figure 1 shows that, especially for the Stacked Embeddings architecture (Torres Aguilar, 2022), the performance in all settings is comparable and closer to the supervised performance of the architecture, remaining above 80% of the supervised F1-Score. This indicates that CuPUL enables the supervised NER architecture to generalize well to weak supervision.

The relatively low capacity of Stacked Embeddings, which is based on a BiLSTM-CRF instead of a transformer, may help in this, as it is, in principle, less prone to overfitting and easier to generalize.

## 7 Conclusion and Future Work

This study shows that weak supervision can effectively be used to train domain-specific models to achieve NER performance levels, on historical texts with little annotations, that are close to fully supervised methods. This is particularly important for the historical domain, where large annotated datasets are often unavailable.

Further assessment of the robustness of this framework and improvement of its performance are the two main areas of focus for future work. Additional analysis could involve examining the adaptability of the framework under various sources of noise that are characteristic of the domain, such as the use of external historical knowledge bases (Uckelman; Wrisley, 2018) to build annotation dictionar-

ies and its application to historical documents transcribed with OCR (Ehrmann et al., 2020). These investigations will determine the framework's applicability when manual transcription of texts or compilation of entity lists is impractical or unfeasible. Additionally, similar analysis can be conducted for other information extraction tasks, including Entity linking and Relationship extraction. In order to increase the NER performance, a major development of the technique would be the integration of active learning, which offers the dual benefit of assisting the model in resolving difficult instances and exploiting domain expertise from historians and linguists where it is most needed and appropriate.

By reducing the need for large annotated datasets, these techniques support the adaptation of NER tools to a wider range of historical low-resource text collections, enabling more comprehensive analysis of historical documents (Bouillon et al., 2024), with significant implications for the fields of digital humanities and cultural heritage.

## Limitations

The empirical analysis of this work is restricted to manually transcribed medieval French charter collections with two entity types (persons and locations), which limits the generalizability of the findings to other historical languages and periods, document genres, larger sets of entity categories, and OCR-transcribed documents. In order to avoid the quality of knowledge bases impacting the results, weak supervision is performed via randomly sampled dictionaries constructed from the training splits; manually curated or external bases, which could be biased towards some subset of entities, may exhibit different noise profiles and cause a different model behavior. Moreover, a fully annotated, albeit small, validation set is kept for model evaluation and hyperparameter tuning. Finally, the evaluation relies on standard NER metrics and does not include human analysis of downstream impact on digital humanities tasks.

## Ethical Considerations

Historical corpora reflect the social, cultural, and political biases of the periods and institutions that produced them. The medieval charters considered in this work primarily document legal and administrative activities of municipal and ecclesiastical authorities, emphasizing the interests of elites. As a result, NER models trained on these sources could reproduce existing biases about people and places, and the usage of models and techniques designed for modern languages or pretrained on modern texts, which is often necessary due to the inherent low volume of available historical data, could amplify these biases when facing challenges specific to the domain, such as linguistic variation and orthographic inconsistency. Moreover, learning from genre-specific conventions like formulaic expressions may capture patterns that are characteristic of the genre and which could not generalize to other genres and time periods.

The systems proposed in this work are intended as assistive tools for historical research, and their outputs should be interpreted in collaboration with domain experts who can assess biases in both the sources and the models.

## Acknowledgments

## References

Baptiste Blouin, Benoit Favre, Jeremy Auguste, and Christian Henriot. 2021. Transferring modern named entity recognition to the historical domain: How to take the step? In *Proceedings of the Workshop on Natural Language Processing for Digital Humanities*, pages 152–162, NIT Silchar, India. NLP Association of India (NLPAI).

Emanuela Boros, Ahmed Hamdi, Elvys Linhares Pontes, Luis Adrián Cabrera-Diego, Jose G. Moreno, Nicolas Sidere, and Antoine Doucet. 2020. Alleviating digitization errors in named entity recognition for historical documents. In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 431–441, Online. Association for Computational Linguistics.

Pierrette Bouillon, Christophe Chazalon, Sandra Coram-Mekkey, Gilles Falquet, Johanna Gerlach, Stephane Marchand-Maillet, Laurent Moccozet, Jonathan Mutal, Raphael Rubino, and Marco Sorbi. 2024. RC-num: A semantic and multilingual online edition of

the geneva council registers from 1545 to 1550. In *Proceedings of the 25th Annual Conference of the European Association for Machine Translation (Volume 2)*, pages 21–22, Sheffield, UK. European Association for Machine Translation (EAMT).

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

de Hemptinne, Thérèse and Deploige, Jeroen and Kupper, Jean-Louis and Prevenier, Walter, editor. 2015. *Diplomata Belgica: les sources diplomatiques des Pays-Bas méridionaux au Moyen Âge. The Diplomatic Sources from the Medieval Southern Low Countries*. Commission royale d'Histoire/Koninklijke Commissie voor Geschiedenis.

Maud Ehrmann, Ahmed Hamdi, Elvys Linhares Pontes, Matteo Romanello, and Antoine Doucet. 2023. Named entity recognition and classification in historical documents: A survey. *ACM Comput. Surv.*, 56(2).

Maud Ehrmann, Matteo Romanello, Alex Flückiger, and Simon Clematide. 2020. Overview of clef hipe 2020: Named entity recognition and linking on historical newspapers. In *Experimental IR Meets Multilinguality, Multimodality, and Interaction: 11th International Conference of the CLEF Association, CLEF 2020, Thessaloniki, Greece, September 22–25, 2020, Proceedings*, page 288–310, Berlin, Heidelberg. Springer-Verlag.

Zheng Fang, Yanan Cao, Tai Li, Ruipeng Jia, Fang Fang, Yanmin Shang, and Yuhai Lu. 2021. TEBNER: Domain specific named entity recognition with type expanded boundary-aware network. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 198–207, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Claire Grover, Sharon Givon, Richard Tobin, and Julian Ball. 2008. Named entity recognition for digitised historical texts. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco. European Language Resources Association (ELRA).

Torsten Hiltmann, Martin Dröge, Nicole Dresselhaus, Till Grallert, Melanie Althage, Paul Bayer, Sophie Eckenstaler, Koray Mendi, Jascha Marijn Schmitz, Philipp Schneider, Wiebke Sczeponik, and Anica Skibba. 2025. NER4all or Context is All You Need: Using LLMs for low-effort, high-performance NER on historical texts. A humanities informed approach. *Preprint*, arXiv:2502.04351.

Imed Keraghel, Stanislas Morbieu, and Mohamed Nadif. 2024. Recent advances in named entity recognition: A comprehensive survey and comparative study. *Preprint*, arXiv:2401.10825.

Kimmo Kettunen, Eetu Mäkelä, Teemu Ruokolainen, Juha Kuokkala, and Laura Löfberg. 2017. Old content and modern tools - searching named entities in a finnish ocred historical newspaper collection 1771-1910. *Digital Humanities Quarterly*, 11.

Ivano Lauriola, Fabio Aiolli, Alberto Lavelli, and Fabio Rinaldi. 2021. Learning adaptive representations for entity recognition in the biomedical domain. *Journal of Biomedical Semantics*, 12(1):10.

Yuepei Li, Kang Zhou, Qiao Qiao, Qing Wang, and Qi Li. 2025. Re-examine distantly supervised NER: A new benchmark and a simple approach. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 10940–10959, Abu Dhabi, UAE. Association for Computational Linguistics.

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, page 1054–1064, New York, NY, USA. Association for Computing Machinery.

Eliana Magnani. 2020. *Des chartae au Corpus: La plateforme des CBMA - Chartae/Corpus Burgundiae Medii Aevi*, volume 27 of *Atelier de Recherche Sur Les Textes Médiévaux*, pages 57–67. Brepols Publishers.

Hiroki Nakayama. 2018. seqeval: A python framework for sequence labeling evaluation. Software available from https://github.com/chakki-works/seqeval.

Vit Novotny, Kristina Luger, Michal Štefánik, Tereza Vrabcova, and Ales Horak. 2023. People and places of historical Europe: Bootstrapping annotation pipeline and a new corpus of named entities in late medieval texts. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 14104–14113, Toronto, Canada. Association for Computational Linguistics.

Minlong Peng, Xiaoyu Xing, Qi Zhang, Jinlan Fu, and Xuanjing Huang. 2019. Distantly supervised named entity recognition using positive-unlabeled learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2409–2419, Florence, Italy. Association for Computational Linguistics.

Cristian Santini, Laura Melosi, and Emanuele Frontoni. 2025. Named entity recognition in historical italian: The case of giacomo leopardi's zibaldone. *Preprint*, arXiv:2505.20113.

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.

Shuzheng Si, Zefan Cai, Shuang Zeng, Guoqiang Feng, Jiaxing Lin, and Baobao Chang. 2023. SANTA: Separate strategies for inaccurate and incomplete annotation noise in distantly-supervised named entity recognition. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3883–3896, Toronto, Canada. Association for Computational Linguistics.

Shuzheng Si, Helan Hu, Haozhe Zhao, Shuang Zeng, Kaikai An, Zefan Cai, and Baobao Chang. 2024. Improving the robustness of distantly-supervised named entity recognition via uncertainty-aware teacher learning and student-student collaborative learning. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 5533–5546, Bangkok, Thailand. Association for Computational Linguistics.

Dominique Stutzmann, Sergio Torres Aguilar, and Paul Chaffenet. 2021. HOME-Alcar: Aligned and Annotated Cartularies. Type: dataset.

Sergio Torres Aguilar. 2022. Multilingual named entity recognition for medieval charters using stacked embeddings and bert-based models. In *Proceedings of the Second Workshop on Language Technologies for Historical and Ancient Languages*, pages 119–128, Marseille, France. European Language Resources Association.

Sergio Torres Aguilar, Xavier Tannier, and Pierre Chastang. 2016. Named entity recognition applied on a data base of Medieval Latin charters. The case of chartae burgundiae. In *3rd International Workshop on Computational History (HistoInformatics 2016)*, Krakow, Poland.

Sara L Uckelman. DRAFT: Names in the 1292 census of Paris.

David Joseph Wrisley. 2018. The literary geographies of christine de pizan. MLA.

Shuhui Wu, Yongliang Shen, Zeqi Tan, Wenqi Ren, Jietian Guo, Shiliang Pu, and Weiming Lu. 2023. MProto: Multi-prototype network with denoised optimal transport for distantly supervised named entity recognition. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2361–2374, Singapore. Association for Computational Linguistics.

Vikas Yadav and Steven Bethard. 2018. A survey on recent advances in named entity recognition from deep learning models. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2145–2158, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Shibingfeng Zhang and Giovanni Colavizza. 2025. Named entity recognition of historical text via large language model. *Preprint*, arXiv:2508.18090.

Xinghua Zhang, Bowen Yu, Tingwen Liu, Zhenyu Zhang, Jiawei Sheng, Xue Mengge, and Hongbo Xu. 2021. Improving distantly-supervised named entity recognition with self-collaborative denoising learning. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1518–1529, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Kang Zhou, Yuepei Li, and Qi Li. 2022. Distantly supervised named entity recognition via confidence-based multi-class positive and unlabeled learning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7198–7211, Dublin, Ireland. Association for Computational Linguistics.

# A  Detailed Results

Tables 3 to 5 report the micro-averaged F1-Score, Precision, and Recall obtained by the different techniques, averaged on the test split of each source, in the cases of fully supervised NER and varying distant supervision settings. These measures are used to draw Figure 1.

Tables 6 to 8 report the micro-averaged F1-Score, Precision, and Recall obtained by the different techniques on the test split of each source, in the cases of fully supervised NER and varying distant supervision settings.

# B  Hyperparameter tuning

For RQ1, we used the default hyperparameters for every technique, computing prior class probabilities from the validation set when required. For RQ2, according to Li et al. (2025) and using grid searches, we made a first tuning step for the voters' hyperparameters, and a second step for the curriculum training hyperparameters, as we used CuPUL with Stacked embeddings (Torres Aguilar, 2022) and XLM-RoBERTa (Conneau et al., 2020) voters.

## B.1  Voters' hyperparameter tuning

The tuned hyperparameters for the voters are the train epochs, drop negative, loss type, and m. The learning rate is set to $1e-3$, and other hyperparameters are left at their default values. Tables 9 to 11 show the validation performances of Stacked Embeddings and XLM-RoBERTa voters with every set of hyperparameters. We select the set of hyperparameters based on the mean and standard deviation of the performances and on the performance with

| Sample (%) | 10 | 20 | 40 | 60 | 80 | 100 | Fully | MEAN |
|---|---|---|---|---|---|---|---|---|
| Conf-MPU | $38 \pm 27$ | $24 \pm 33$ | $22 \pm 31$ | $36 \pm 27$ | $29 \pm 22$ | $17 \pm 24$ | $43 \pm 31$ | $30 \pm 8$ |
| CuPUL | $\mathbf{67} \pm 1$ | $\mathbf{68} \pm 0$ | $\mathbf{67} \pm 3$ | $68 \pm 5$ | $67 \pm 5$ | $66 \pm 5$ | $84 \pm 7$ | $\mathbf{69} \pm 6$ |
| BOND | $32 \pm 18$ | $57 \pm 13$ | $\mathbf{67} \pm 5$ | $\mathbf{74} \pm 5$ | $\mathbf{73} \pm 5$ | $\mathbf{77} \pm 11$ | $\mathbf{90} \pm 4$ | $67 \pm 16$ |
| SCDL | $0 \pm 0$ | $0 \pm 0$ | $9 \pm 12$ | $28 \pm 39$ | $28 \pm 39$ | $28 \pm 38$ | $29 \pm 40$ | $18 \pm 12$ |
| CENSOR | $4 \pm 5$ | $10 \pm 7$ | $33 \pm 32$ | $50 \pm 36$ | $52 \pm 37$ | $53 \pm 38$ | $60 \pm 42$ | $37 \pm 19$ |
| SANTA | $46 \pm 0$ | $54 \pm 3$ | $63 \pm 3$ | $68 \pm 6$ | $72 \pm 4$ | $76 \pm 5$ | $\mathbf{90} \pm 4$ | $67 \pm 13$ |
| MProto | $45 \pm 6$ | $53 \pm 3$ | $56 \pm 9$ | $60 \pm 10$ | $52 \pm 15$ | $59 \pm 8$ | $64 \pm 18$ | $56 \pm 5$ |
| StEmb | $17 \pm 11$ | $38 \pm 20$ | $54 \pm 6$ | $73 \pm 8$ | $80 \pm 3$ | $85 \pm 7$ | $95 \pm 1$ | $63 \pm 24$ |
| Matching | $16 \pm 9$ | $24 \pm 7$ | $33 \pm 4$ | $41 \pm 2$ | $46 \pm 3$ | $51 \pm 5$ | | |
| CuPUL+XLM | $68 \pm 6$ | $78 \pm 3$ | $78 \pm 2$ | $76 \pm 5$ | $\mathbf{77} \pm 5$ | $79 \pm 7$ | $\mathbf{92} \pm 3$ | $78 \pm 6$ |
| CuPUL+StEmb | $\mathbf{78} \pm 3$ | $\mathbf{82} \pm 3$ | $\mathbf{79} \pm 1$ | $\mathbf{80} \pm 4$ | $77 \pm 6$ | $76 \pm 7$ | $88 \pm 4$ | $\mathbf{80} \pm 3$ |

Table 3: Test results in terms of F1-Score of Fully supervised NER and Weakly supervised NER varying the dictionary sample size. Measures are micro-averaged among classes and then, for each setting, averaged among datasets and reported in percentage with their standard deviation. The MEAN column reports the performance averaged across settings with its standard deviation. The best results among DS-NER techniques for each setting are reported in bold. For reference, the Stacked Embeddings (StEmb) supervised architecture and the pure dictionary matching are reported, too. The bottom rows report performances of the CuPUL technique paired with XLM-RoBERTa (XLM) and Stacked Embeddings, respectively, as underlying classification model.

low supervision settings. The selected hyperparameters are:

- Stacked Embeddings voter hyperparameters:
    - train epochs (E) = 15, drop negative (Neg) = 0.1, loss type = MPN, m = 20

- XLM-RoBERTa voter hyperparameters:
    - train epochs (E) = 5, drop negative (Neg) = 0.5, loss type = MPN-CE, m = 20

## B.2 Curriculum training hyperparameter tuning

Tables 12 and 13 show the validation performances of the CuPUL curriculum train stage with Stacked Embeddings and XLM-RoBERTa, respectively, with every set of hyperparameters. We select the set of hyperparameters based on the mean and standard deviation of the performances and on the performance with low supervision settings. The selected hyperparameters are:

- CuPUL with Stacked Embeddings:
    - train epochs (E) = 20, loss type = MPN

- CuPUL with XLM-RoBERTa:
    - train epochs (E) = 10, train sub-epochs (sE) = 1, loss type = Conf-MPU-CE, learning rate (LR): 1e−5

## C Online Resources

The source code is available online at https://github.com/msorbi/hwsner

| Sample (%) | 10 | 20 | 40 | 60 | 80 | 100 | Fully | MEAN |
|---|---|---|---|---|---|---|---|---|
| Conf-MPU | $68 \pm 23$ | $55 \pm 41$ | $52 \pm 41$ | $67 \pm 24$ | $59 \pm 29$ | $46 \pm 41$ | $71 \pm 21$ | $60 \pm 8$ |
| CuPUL | $66 \pm 2$ | $62 \pm 2$ | $58 \pm 4$ | $58 \pm 6$ | $55 \pm 6$ | $54 \pm 6$ | $81 \pm 7$ | $62 \pm 8$ |
| BOND | $\mathbf{76} \pm 18$ | $\mathbf{70} \pm 8$ | $66 \pm 5$ | $\mathbf{69} \pm 6$ | $\mathbf{66} \pm 10$ | $\mathbf{68} \pm 16$ | $\mathbf{88} \pm 6$ | $\mathbf{72} \pm 7$ |
| SCDL | $0 \pm 0$ | $0 \pm 0$ | $12 \pm 17$ | $26 \pm 37$ | $26 \pm 36$ | $59 \pm 42$ | $28 \pm 39$ | $22 \pm 18$ |
| CENSOR | $73 \pm 33$ | $64 \pm 28$ | $\mathbf{70} \pm 24$ | $47 \pm 34$ | $48 \pm 34$ | $48 \pm 34$ | $58 \pm 41$ | $58 \pm 10$ |
| SANTA | $64 \pm 10$ | $59 \pm 8$ | $57 \pm 5$ | $59 \pm 7$ | $62 \pm 5$ | $65 \pm 8$ | $87 \pm 5$ | $65 \pm 9$ |
| MProto | $45 \pm 7$ | $51 \pm 7$ | $49 \pm 7$ | $52 \pm 10$ | $42 \pm 14$ | $49 \pm 8$ | $56 \pm 21$ | $49 \pm 4$ |
| StEmb | $77 \pm 4$ | $70 \pm 9$ | $66 \pm 1$ | $76 \pm 8$ | $76 \pm 7$ | $78 \pm 11$ | $95 \pm 1$ | $77 \pm 8$ |
| Matching | $65 \pm 17$ | $60 \pm 9$ | $52 \pm 7$ | $52 \pm 10$ | $52 \pm 9$ | $52 \pm 9$ | | |
| CuPUL+XLM | $74 \pm 3$ | $77 \pm 4$ | $71 \pm 4$ | $67 \pm 6$ | $68 \pm 7$ | $\mathbf{69} \pm 10$ | $92 \pm 3$ | $\mathbf{74} \pm 7$ |
| CuPUL+StEmb | $\mathbf{76} \pm 2$ | $\mathbf{79} \pm 4$ | $\mathbf{72} \pm 3$ | $\mathbf{72} \pm 6$ | $\mathbf{69} \pm 8$ | $67 \pm 10$ | $86 \pm 5$ | $\mathbf{74} \pm 5$ |

Table 4: Test results in terms of Precision of Fully supervised NER and Weakly supervised NER varying the dictionary sample size. Measures are micro-averaged among classes and then, for each setting, averaged among datasets and reported in percentage with their standard deviation. The MEAN column reports the performance averaged across settings with its standard deviation. The best results among DS-NER techniques for each setting are reported in bold. For reference, the Stacked Embeddings (StEmb) supervised architecture and the pure dictionary matching are reported, too. The bottom rows report performances of the CuPUL technique paired with XLM-RoBERTa (XLM) and Stacked Embeddings, respectively, as underlying classification model.

| Sample (%) | 10 | 20 | 40 | 60 | 80 | 100 | Fully | MEAN |
|---|---|---|---|---|---|---|---|---|
| Conf-MPU | $42 \pm 30$ | $26 \pm 37$ | $28 \pm 39$ | $41 \pm 34$ | $38 \pm 32$ | $27 \pm 39$ | $48 \pm 35$ | $36 \pm 7$ |
| CuPUL | $\mathbf{68} \pm 3$ | $\mathbf{74} \pm 2$ | $\mathbf{79} \pm 6$ | $\mathbf{83} \pm 6$ | $84 \pm 5$ | $85 \pm 5$ | $88 \pm 6$ | $\mathbf{80} \pm 6$ |
| BOND | $23 \pm 16$ | $50 \pm 19$ | $71 \pm 14$ | $81 \pm 8$ | $84 \pm 5$ | $92 \pm 1$ | $93 \pm 2$ | $71 \pm 22$ |
| SCDL | $1 \pm 1$ | $1 \pm 1$ | $8 \pm 9$ | $30 \pm 41$ | $30 \pm 41$ | $31 \pm 42$ | $31 \pm 42$ | $19 \pm 13$ |
| CENSOR | $2 \pm 3$ | $6 \pm 4$ | $34 \pm 38$ | $53 \pm 39$ | $56 \pm 40$ | $59 \pm 42$ | $61 \pm 43$ | $39 \pm 22$ |
| SANTA | $37 \pm 4$ | $51 \pm 4$ | $69 \pm 3$ | $80 \pm 6$ | $\mathbf{87} \pm 4$ | $\mathbf{93} \pm 2$ | $\mathbf{94} \pm 3$ | $73 \pm 19$ |
| MProto | $45 \pm 5$ | $57 \pm 5$ | $66 \pm 14$ | $70 \pm 12$ | $68 \pm 15$ | $74 \pm 10$ | $78 \pm 9$ | $66 \pm 10$ |
| StEmb | $10 \pm 7$ | $29 \pm 19$ | $46 \pm 8$ | $70 \pm 8$ | $84 \pm 3$ | $95 \pm 1$ | $95 \pm 1$ | $61 \pm 29$ |
| Matching | $9 \pm 6$ | $15 \pm 5$ | $25 \pm 4$ | $35 \pm 2$ | $42 \pm 3$ | $52 \pm 2$ | | |
| CuPUL+XLM | $64 \pm 9$ | $80 \pm 5$ | $85 \pm 3$ | $88 \pm 3$ | $\mathbf{91} \pm 2$ | $\mathbf{92} \pm 2$ | $\mathbf{93} \pm 2$ | $85 \pm 9$ |
| CuPUL+StEmb | $\mathbf{80} \pm 5$ | $\mathbf{84} \pm 2$ | $\mathbf{88} \pm 1$ | $\mathbf{89} \pm 1$ | $89 \pm 1$ | $89 \pm 2$ | $91 \pm 3$ | $\mathbf{87} \pm 4$ |

Table 5: Test results in terms of Recall of Fully supervised NER and Weakly supervised NER varying the dictionary sample size. Measures are micro-averaged among classes and then, for each setting, averaged among datasets and reported in percentage with their standard deviation. The MEAN column reports the performance averaged across settings with its standard deviation. The best results among DS-NER techniques for each setting are reported in bold. For reference, the Stacked Embeddings (StEmb) supervised architecture and the pure dictionary matching are reported, too. The bottom rows report performances of the CuPUL technique paired with XLM-RoBERTa (XLM) and Stacked Embeddings, respectively, as underlying classification model.

| Sample (%) | 10 | 20 | 40 | 60 | 80 | 100 | Fully | MEAN |
|---|---|---|---|---|---|---|---|---|
| Conf-MPU | 55 | 0 | 0 | 43 | 37 | 0 | 71 | 30 ± 26 |
| CuPUL | **66** | **67** | **62** | 63 | 63 | 64 | 75 | **66 ± 4** |
| BOND | 29 | 43 | 60 | **70** | **75** | **86** | **90** | 65 ± 19 |
| SCDL | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 ± 0 |
| CENSOR | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 ± 0 |
| SANTA | 46 | 53 | 59 | 59 | 68 | 79 | 85 | 64 ± 12 |
| MProto | 37 | 49 | 47 | 46 | 31 | 47 | 39 | 42 ± 6 |
| StEmb | 32 | 66 | 59 | 83 | 83 | 92 | 93 | 73 ± 19 |
| Matching | 27 | 33 | 38 | 41 | 44 | 51 | | |
| CuPUL+XLM | 77 | **82** | **80** | 75 | 77 | **82** | **89** | 80 ± 4 |
| CuPUL+StEmb | **80** | 81 | **80** | 84 | 80 | 80 | 83 | **81 ± 1** |
| | | | | | | | | |
| Conf-MPU | 0 | 71 | 67 | 64 | 51 | 51 | 56 | 52 ± 21 |
| CuPUL | **68** | 68 | 69 | 65 | 63 | 61 | 89 | **69 ± 8** |
| BOND | 55 | **74** | 72 | 73 | 66 | 61 | 86 | **69 ± 9** |
| SCDL | 0 | 0 | 26 | 82 | 82 | **82** | 86 | 51 ± 35 |
| CENSOR | 1 | 12 | **75** | **85** | **85** | 76 | 91 | 61 ± 33 |
| SANTA | 45 | 59 | 64 | 71 | 72 | 68 | **94** | 68 ± 13 |
| MProto | 51 | 57 | 69 | 64 | 62 | 64 | 83 | 64 ± 9 |
| StEmb | 6 | 22 | 57 | 66 | 76 | 76 | 96 | 57 ± 28 |
| Matching | 5 | 17 | 33 | 39 | 43 | 45 | | |
| CuPUL+XLM | 65 | 74 | 75 | 72 | **71** | **69** | 93 | 74 ± 8 |
| CuPUL+StEmb | **80** | 78 | 78 | 74 | 69 | 67 | 91 | **77 ± 7** |
| | | | | | | | | |
| Conf-MPU | 57 | 0 | 0 | 0 | 0 | 0 | 0 | 8 ± 19 |
| CuPUL | **66** | **68** | 69 | 75 | 73 | 73 | 88 | **73 ± 6** |
| BOND | 12 | 53 | **70** | **81** | 78 | 85 | 95 | 68 ± 24 |
| SCDL | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 ± 0 |
| CENSOR | 11 | 17 | 22 | 64 | 69 | 83 | 87 | 51 ± 28 |
| SANTA | 46 | 51 | 65 | 74 | **78** | 80 | 91 | 69 ± 14 |
| MProto | 47 | 53 | 53 | 70 | 63 | 66 | 69 | 60 ± 8 |
| StEmb | 14 | 27 | 46 | 70 | 81 | 88 | 95 | 60 ± 27 |
| Matching | 15 | 23 | 29 | 44 | 50 | 57 | | |
| CuPUL+XLM | 64 | 78 | 78 | **83** | **84** | **85** | 95 | 81 ± 8 |
| CuPUL+StEmb | **73** | **86** | 80 | 82 | 83 | 83 | 91 | **82 ± 5** |

Table 6: Test results in terms of F1-Score of Fully supervised NER and Weakly supervised NER varying the dictionary sample size on each of the three datasets (CBMA, CDBE, HOME), respectively. Measures are micro-averaged among classes. The MEAN column reports the performance averaged across settings with its standard deviation. The best results among DS-NER techniques for each setting are reported in bold. For reference, the Stacked Embeddings (StEmb) supervised architecture and the pure dictionary matching are reported, too. The bottom rows report performances of the CuPUL technique paired with XLM-RoBERTa (XLM) and Stacked Embeddings, respectively, as underlying classification model.

| Sample (%) | 10 | 20 | 40 | 60 | 80 | 100 | Fully | MEAN |
|---|---|---|---|---|---|---|---|---|
| Conf-MPU | 52 | 0 | 0 | 48 | 39 | 0 | 63 | $29 \pm 24$ |
| CuPUL | 64 | 62 | 54 | 55 | 53 | 54 | 70 | $59 \pm 6$ |
| BOND | 57 | 59 | 65 | **69** | **73** | **80** | 89 | **70** $\pm 10$ |
| SCDL | 1 | 1 | 1 | 1 | 1 | 1 | 1 | $1 \pm 0$ |
| CENSOR | **100** | **100** | **100** | 0 | 0 | 0 | 0 | $43 \pm 46$ |
| SANTA | 51 | 50 | 52 | 50 | 58 | 70 | 80 | $59 \pm 10$ |
| MProto | 35 | 42 | 41 | 40 | 23 | 39 | 28 | $35 \pm 6$ |
| StEmb | 74 | 80 | 68 | 85 | 83 | 89 | 93 | $82 \pm 8$ |
| Matching | 74 | 66 | 58 | 52 | 52 | 53 | | |
| CuPUL+XLM | **77** | **78** | 73 | 66 | 68 | **74** | **87** | $75 \pm 6$ |
| CuPUL+StEmb | 76 | 76 | **75** | **78** | **72** | 73 | 79 | **76** $\pm 2$ |
| | | | | | | | | |
| Conf-MPU | **100** | 65 | 56 | 52 | 38 | 37 | 51 | $57 \pm 18$ |
| CuPUL | 64 | 60 | 57 | 52 | 49 | 46 | 86 | $59 \pm 12$ |
| BOND | 73 | **72** | 60 | 62 | 52 | 46 | 80 | $63 \pm 11$ |
| SCDL | 0 | 0 | 37 | **78** | 77 | **75** | 83 | $50 \pm 32$ |
| CENSOR | 28 | 32 | **67** | **78** | **79** | 67 | 90 | $63 \pm 21$ |
| SANTA | 74 | 68 | 57 | 60 | 59 | 53 | **92** | **66** $\pm 12$ |
| MProto | 49 | 54 | 58 | 52 | 49 | 51 | 77 | $56 \pm 9$ |
| StEmb | 74 | 57 | 66 | 66 | 66 | 64 | 96 | $70 \pm 11$ |
| Matching | 41 | 48 | 42 | 41 | 41 | 40 | | |
| CuPUL+XLM | 70 | 72 | 66 | 61 | **59** | **55** | 93 | $68 \pm 11$ |
| CuPUL+StEmb | **79** | **76** | **69** | **63** | 57 | 54 | 90 | **70** $\pm 11$ |
| | | | | | | | | |
| Conf-MPU | 50 | **100** | **100** | **100** | **100** | **100** | **100** | **93** $\pm 16$ |
| CuPUL | 69 | 65 | 63 | 66 | 63 | 62 | 85 | $68 \pm 7$ |
| BOND | **100** | 78 | 72 | 77 | 71 | 78 | 95 | $82 \pm 10$ |
| SCDL | 0 | 0 | 0 | 0 | 0 | **100** | 0 | $14 \pm 33$ |
| CENSOR | 93 | 59 | 42 | 63 | 64 | 78 | 85 | $69 \pm 15$ |
| SANTA | 68 | 57 | 63 | 67 | 70 | 70 | 88 | $69 \pm 8$ |
| MProto | 51 | 57 | 48 | 65 | 56 | 58 | 62 | $57 \pm 5$ |
| StEmb | 83 | 72 | 66 | 78 | 80 | 82 | 96 | $79 \pm 8$ |
| Matching | 80 | 67 | 55 | 65 | 62 | 62 | | |
| CuPUL+XLM | **76** | 82 | **76** | 75 | **76** | 78 | 94 | **80** $\pm 6$ |
| CuPUL+StEmb | 73 | **85** | 74 | **76** | **76** | 75 | 88 | $78 \pm 5$ |

Table 7: Test results in terms of Precision of Fully supervised NER and Weakly supervised NER varying the dictionary sample size on each of the three datasets (CBMA, CDBE, HOME), respectively. Measures are micro-averaged among classes. The MEAN column reports the performance averaged across settings with its standard deviation. The best results among DS-NER techniques for each setting are reported in bold. For reference, the Stacked Embeddings (StEmb) supervised architecture and the pure dictionary matching are reported, too. The bottom rows report performances of the CuPUL technique paired with XLM-RoBERTa (XLM) and Stacked Embeddings, respectively, as underlying classification model.

| Sample (%) | 10 | 20 | 40 | 60 | 80 | 100 | Fully | MEAN |
|---|---|---|---|---|---|---|---|---|
| Conf-MPU | 59 | 0 | 0 | 39 | 35 | 0 | 81 | $31 \pm 28$ |
| CuPUL | **68** | **74** | **73** | **74** | 77 | 78 | 80 | **75 $\pm$ 3** |
| BOND | 20 | 33 | 55 | 70 | 77 | **92** | 92 | $63 \pm 24$ |
| SCDL | 2 | 2 | 2 | 2 | 3 | 3 | 3 | $2 \pm 0$ |
| CENSOR | 0 | 0 | 0 | 0 | 0 | 0 | 0 | $0 \pm 0$ |
| SANTA | 42 | 55 | 67 | 71 | **82** | 90 | 90 | $71 \pm 15$ |
| MProto | 40 | 60 | 53 | 54 | 49 | 61 | 67 | $55 \pm 8$ |
| StEmb | 20 | 56 | 52 | 82 | 83 | 95 | 93 | $69 \pm 24$ |
| Matching | 17 | 22 | 28 | 34 | 39 | 49 | | |
| CuPUL+XLM | 76 | **86** | **87** | 85 | **88** | **92** | **91** | $86 \pm 4$ |
| CuPUL+StEmb | **84** | **86** | 86 | **90** | **88** | 88 | 87 | **87 $\pm$ 2** |
| | | | | | | | | |
| Conf-MPU | 0 | **78** | 83 | 84 | 79 | 82 | 63 | $67 \pm 26$ |
| CuPUL | **72** | **78** | 87 | 86 | 88 | 88 | 92 | **84 $\pm$ 6** |
| BOND | 45 | 77 | **90** | 88 | 89 | 91 | 92 | $82 \pm 15$ |
| SCDL | 0 | 0 | 21 | 87 | 89 | 90 | 90 | $54 \pm 38$ |
| CENSOR | 0 | 8 | 87 | **92** | **93** | 89 | 93 | $66 \pm 37$ |
| SANTA | 33 | 52 | 73 | 86 | 92 | **95** | **96** | $75 \pm 21$ |
| MProto | 53 | 61 | 85 | 81 | 84 | 85 | 90 | $77 \pm 12$ |
| StEmb | 3 | 13 | 50 | 66 | 88 | 94 | 96 | $59 \pm 33$ |
| Matching | 3 | 10 | 27 | 37 | 45 | 53 | | |
| CuPUL+XLM | 61 | 77 | 87 | 87 | **90** | **90** | **94** | $84 \pm 10$ |
| CuPUL+StEmb | **82** | 81 | **90** | **89** | 89 | 88 | 93 | **87 $\pm$ 4** |
| | | | | | | | | |
| Conf-MPU | **66** | 0 | 0 | 0 | 0 | 0 | 0 | $9 \pm 22$ |
| CuPUL | 64 | **72** | **77** | **87** | 87 | 89 | 91 | **81 $\pm$ 9** |
| BOND | 6 | 40 | 68 | 85 | 86 | **94** | 95 | $68 \pm 29$ |
| SCDL | 0 | 0 | 0 | 0 | 0 | 0 | 0 | $0 \pm 0$ |
| CENSOR | 6 | 10 | 15 | 66 | 76 | 90 | 90 | $50 \pm 33$ |
| SANTA | 34 | 46 | 67 | 82 | **87** | 93 | 94 | $72 \pm 20$ |
| MProto | 43 | 50 | 60 | 75 | 72 | 77 | 77 | $65 \pm 12$ |
| StEmb | 8 | 17 | 35 | 63 | 81 | 94 | 95 | $56 \pm 32$ |
| Matching | 8 | 14 | 19 | 33 | 42 | 53 | | |
| CuPUL+XLM | 55 | 75 | 81 | **92** | **94** | **95** | **96** | $84 \pm 13$ |
| CuPUL+StEmb | **73** | **86** | **87** | 89 | 92 | 92 | 94 | **87 $\pm$ 6** |

Table 8: Test results in terms of Recall of Fully supervised NER and Weakly supervised NER varying the dictionary sample size on each of the three datasets (CBMA, CDBE, HOME), respectively. Measures are micro-averaged among classes. The MEAN column reports the performance averaged across settings with its standard deviation. The best results among DS-NER techniques for each setting are reported in bold. For reference, the Stacked Embeddings (StEmb) supervised architecture and the pure dictionary matching are reported, too. The bottom rows report performances of the CuPUL technique paired with XLM-RoBERTa (XLM) and Stacked Embeddings, respectively, as underlying classification model.

| E | Neg | Loss | m | 10 | 20 | 40 | 60 | 80 | 100 | Fully | MEAN |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 10 | 0.1 | MPN | 10 | $53 \pm 37$ | $54 \pm 38$ | $55 \pm 39$ | $54 \pm 38$ | $83 \pm 5$ | $82 \pm 6$ | $92 \pm 2$ | $68 \pm 15$ |
| 10 | 0.1 | MPN | 20 | $76 \pm 5$ | $\mathbf{82} \pm 5$ | $81 \pm 3$ | $81 \pm 3$ | $80 \pm 6$ | $80 \pm 7$ | $90 \pm 2$ | $81 \pm 4$ |
| 10 | 0.1 | MPN-CE | 10 | $75 \pm 1$ | $79 \pm 2$ | $83 \pm 3$ | $83 \pm 3$ | $84 \pm 6$ | $83 \pm 7$ | $93 \pm 2$ | $\mathbf{83} \pm 5$ |
| 10 | 0.1 | MPN-CE | 20 | $73 \pm 5$ | $79 \pm 3$ | $81 \pm 2$ | $82 \pm 3$ | $81 \pm 6$ | $80 \pm 8$ | $92 \pm 1$ | $81 \pm 5$ |
| 10 | 0.3 | MPN | 10 | $54 \pm 38$ | $53 \pm 37$ | $55 \pm 39$ | $54 \pm 38$ | $83 \pm 5$ | $82 \pm 7$ | $92 \pm 3$ | $67 \pm 15$ |
| 10 | 0.3 | MPN | 20 | $74 \pm 9$ | $81 \pm 6$ | $81 \pm 4$ | $80 \pm 3$ | $80 \pm 5$ | $78 \pm 7$ | $90 \pm 3$ | $80 \pm 4$ |
| 10 | 0.3 | MPN-CE | 10 | $75 \pm 2$ | $79 \pm 1$ | $83 \pm 3$ | $83 \pm 3$ | $83 \pm 5$ | $83 \pm 7$ | $93 \pm 1$ | $\mathbf{83} \pm 5$ |
| 10 | 0.3 | MPN-CE | 20 | $73 \pm 4$ | $80 \pm 2$ | $82 \pm 2$ | $82 \pm 3$ | $81 \pm 6$ | $81 \pm 8$ | $92 \pm 1$ | $81 \pm 5$ |
| 15 | 0.1 | MPN | 10 | $53 \pm 37$ | $55 \pm 39$ | $\mathbf{87} \pm 2$ | $\mathbf{86} \pm 3$ | $85 \pm 5$ | $83 \pm 6$ | $92 \pm 2$ | $77 \pm 14$ |
| **15** | **0.1** | **MPN** | **20** | $\mathbf{80} \pm 2$ | $81 \pm 7$ | $84 \pm 2$ | $81 \pm 2$ | $79 \pm 6$ | $80 \pm 8$ | $91 \pm 1$ | $82 \pm 3$ |
| 15 | 0.1 | MPN-CE | 10 | $69 \pm 5$ | $75 \pm 4$ | $81 \pm 3$ | $85 \pm 5$ | $84 \pm 7$ | $84 \pm 8$ | $\mathbf{94} \pm 0$ | $82 \pm 7$ |
| 15 | 0.1 | MPN-CE | 20 | $71 \pm 5$ | $75 \pm 6$ | $79 \pm 0$ | $83 \pm 4$ | $82 \pm 6$ | $82 \pm 8$ | $92 \pm 1$ | $81 \pm 6$ |
| 15 | 0.3 | MPN | 10 | $51 \pm 36$ | $53 \pm 37$ | $\mathbf{87} \pm 2$ | $85 \pm 3$ | $85 \pm 5$ | $\mathbf{85} \pm 7$ | $92 \pm 2$ | $77 \pm 15$ |
| 15 | 0.3 | MPN | 20 | $78 \pm 2$ | $81 \pm 6$ | $83 \pm 1$ | $82 \pm 3$ | $79 \pm 6$ | $81 \pm 7$ | $90 \pm 1$ | $82 \pm 3$ |
| 15 | 0.3 | MPN-CE | 10 | $69 \pm 4$ | $75 \pm 2$ | $81 \pm 3$ | $85 \pm 5$ | $84 \pm 7$ | $84 \pm 8$ | $\mathbf{94} \pm 1$ | $82 \pm 7$ |
| 15 | 0.3 | MPN-CE | 20 | $70 \pm 5$ | $75 \pm 6$ | $81 \pm 2$ | $83 \pm 4$ | $82 \pm 6$ | $81 \pm 9$ | $92 \pm 1$ | $81 \pm 6$ |
| 20 | 0.1 | MPN | 10 | $52 \pm 37$ | $53 \pm 38$ | $84 \pm 3$ | $\mathbf{86} \pm 3$ | $\mathbf{86} \pm 5$ | $84 \pm 6$ | $93 \pm 1$ | $77 \pm 15$ |
| 20 | 0.1 | MPN | 20 | $77 \pm 2$ | $81 \pm 7$ | $84 \pm 2$ | $83 \pm 3$ | $82 \pm 6$ | $82 \pm 8$ | $91 \pm 1$ | $\mathbf{83} \pm 3$ |
| 20 | 0.1 | MPN-CE | 10 | $67 \pm 6$ | $71 \pm 6$ | $80 \pm 3$ | $84 \pm 5$ | $84 \pm 6$ | $\mathbf{85} \pm 9$ | $\mathbf{94} \pm 1$ | $81 \pm 8$ |
| 20 | 0.1 | MPN-CE | 20 | $68 \pm 7$ | $73 \pm 6$ | $79 \pm 2$ | $83 \pm 5$ | $82 \pm 6$ | $83 \pm 9$ | $\mathbf{94} \pm 1$ | $80 \pm 7$ |
| 20 | 0.3 | MPN | 10 | $51 \pm 36$ | $53 \pm 38$ | $\mathbf{87} \pm 3$ | $\mathbf{86} \pm 3$ | $\mathbf{86} \pm 5$ | $\mathbf{85} \pm 7$ | $92 \pm 1$ | $77 \pm 15$ |
| 20 | 0.3 | MPN | 20 | $79 \pm 0$ | $81 \pm 7$ | $82 \pm 1$ | $83 \pm 4$ | $81 \pm 6$ | $81 \pm 8$ | $91 \pm 2$ | $\mathbf{83} \pm 3$ |
| 20 | 0.3 | MPN-CE | 10 | $67 \pm 5$ | $72 \pm 5$ | $81 \pm 5$ | $84 \pm 6$ | $84 \pm 6$ | $\mathbf{85} \pm 9$ | $\mathbf{94} \pm 1$ | $81 \pm 8$ |
| 20 | 0.3 | MPN-CE | 20 | $68 \pm 7$ | $73 \pm 6$ | $80 \pm 2$ | $83 \pm 5$ | $82 \pm 7$ | $83 \pm 9$ | $\mathbf{94} \pm 1$ | $81 \pm 7$ |

Table 9: Validation results in terms of F1-Score of Stacked Embeddings voters. Measures are micro-averaged among classes and then, for each setting, averaged among datasets and reported in percentage with their standard deviation. The MEAN column reports the performance averaged across settings with its standard deviation. The best results for each setting are reported in bold. Selected hyperparameters are in bold.

| E | Neg | Loss | m | 10 | 20 | 40 | 60 | 80 | 100 | Fully | MEAN |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 10 | 0.1 | MPN | 10 | $40 \pm 8$ | $68 \pm 1$ | $81 \pm 3$ | $\mathbf{87} \pm 4$ | $90 \pm 3$ | $\mathbf{92} \pm 1$ | $92 \pm 3$ | $78 \pm 16$ |
| 10 | 0.1 | MPN | 20 | $39 \pm 4$ | $69 \pm 3$ | $79 \pm 5$ | $84 \pm 7$ | $90 \pm 3$ | $89 \pm 4$ | $92 \pm 4$ | $77 \pm 16$ |
| 10 | 0.1 | MPN-CE | 10 | $25 \pm 5$ | $40 \pm 7$ | $56 \pm 8$ | $77 \pm 5$ | $87 \pm 2$ | $88 \pm 3$ | $94 \pm 3$ | $67 \pm 23$ |
| 10 | 0.1 | MPN-CE | 20 | $32 \pm 11$ | $44 \pm 5$ | $60 \pm 9$ | $77 \pm 2$ | $87 \pm 0$ | $91 \pm 0$ | $93 \pm 4$ | $69 \pm 21$ |
| 10 | 0.1 | MPU | 10 | $69 \pm 11$ | $55 \pm 19$ | $47 \pm 1$ | $45 \pm 2$ | $43 \pm 3$ | $36 \pm 7$ | $42 \pm 4$ | $48 \pm 9$ |
| 10 | 0.1 | MPU | 20 | $65 \pm 17$ | $50 \pm 10$ | $56 \pm 15$ | $52 \pm 14$ | $39 \pm 1$ | $39 \pm 1$ | $48 \pm 8$ | $50 \pm 8$ |
| 10 | 0.1 | MPU-CE | 10 | $64 \pm 9$ | $60 \pm 16$ | $61 \pm 7$ | $54 \pm 13$ | $51 \pm 12$ | $46 \pm 12$ | $40 \pm 16$ | $54 \pm 7$ |
| 10 | 0.1 | MPU-CE | 20 | $58 \pm 13$ | $55 \pm 13$ | $54 \pm 10$ | $45 \pm 15$ | $41 \pm 10$ | $36 \pm 12$ | $43 \pm 9$ | $47 \pm 7$ |
| 10 | 0.3 | MPN | 10 | $31 \pm 7$ | $69 \pm 3$ | $80 \pm 3$ | $86 \pm 6$ | $90 \pm 3$ | $90 \pm 2$ | $93 \pm 2$ | $77 \pm 19$ |
| 10 | 0.3 | MPN | 20 | $35 \pm 9$ | $67 \pm 0$ | $81 \pm 4$ | $86 \pm 5$ | $88 \pm 5$ | $89 \pm 4$ | $91 \pm 3$ | $77 \pm 17$ |
| 10 | 0.3 | MPN-CE | 10 | $27 \pm 7$ | $38 \pm 6$ | $59 \pm 6$ | $77 \pm 4$ | $86 \pm 2$ | $88 \pm 2$ | $94 \pm 2$ | $67 \pm 23$ |
| 10 | 0.3 | MPN-CE | 20 | $32 \pm 11$ | $44 \pm 5$ | $62 \pm 6$ | $78 \pm 3$ | $87 \pm 0$ | $91 \pm 0$ | $94 \pm 3$ | $70 \pm 21$ |
| 10 | 0.3 | MPU | 10 | $\mathbf{70} \pm 11$ | $52 \pm 14$ | $46 \pm 3$ | $47 \pm 0$ | $47 \pm 11$ | $48 \pm 9$ | $44 \pm 3$ | $50 \pm 8$ |
| 10 | 0.3 | MPU | 20 | $63 \pm 20$ | $57 \pm 17$ | $53 \pm 15$ | $55 \pm 13$ | $38 \pm 0$ | $37 \pm 2$ | $41 \pm 1$ | $49 \pm 9$ |
| 10 | 0.3 | MPU-CE | 10 | $65 \pm 9$ | $59 \pm 13$ | $64 \pm 7$ | $55 \pm 14$ | $52 \pm 11$ | $38 \pm 19$ | $51 \pm 9$ | $55 \pm 8$ |
| 10 | 0.3 | MPU-CE | 20 | $58 \pm 11$ | $59 \pm 12$ | $57 \pm 8$ | $45 \pm 15$ | $40 \pm 12$ | $34 \pm 15$ | $56 \pm 14$ | $50 \pm 9$ |
| 10 | 0.5 | MPN | 10 | $36 \pm 11$ | $67 \pm 3$ | $80 \pm 3$ | $85 \pm 2$ | $90 \pm 2$ | $\mathbf{92} \pm 1$ | $93 \pm 3$ | $78 \pm 18$ |
| 10 | 0.5 | MPN | 20 | $40 \pm 8$ | $71 \pm 5$ | $\mathbf{82} \pm 7$ | $85 \pm 6$ | $90 \pm 3$ | $88 \pm 3$ | $93 \pm 3$ | $78 \pm 16$ |
| 10 | 0.5 | MPN-CE | 10 | $26 \pm 6$ | $39 \pm 5$ | $61 \pm 5$ | $78 \pm 3$ | $86 \pm 2$ | $87 \pm 3$ | $94 \pm 2$ | $67 \pm 23$ |
| 10 | 0.5 | MPN-CE | 20 | $33 \pm 11$ | $44 \pm 3$ | $63 \pm 5$ | $77 \pm 2$ | $85 \pm 2$ | $91 \pm 1$ | $\mathbf{95} \pm 2$ | $70 \pm 21$ |
| 10 | 0.5 | MPU | 10 | $69 \pm 12$ | $61 \pm 21$ | $43 \pm 2$ | $44 \pm 1$ | $44 \pm 1$ | $48 \pm 8$ | $45 \pm 2$ | $51 \pm 9$ |
| 10 | 0.5 | MPU | 20 | $63 \pm 17$ | $56 \pm 19$ | $57 \pm 13$ | $44 \pm 8$ | $41 \pm 3$ | $35 \pm 6$ | $45 \pm 7$ | $49 \pm 9$ |
| 10 | 0.5 | MPU-CE | 10 | $63 \pm 9$ | $61 \pm 12$ | $62 \pm 11$ | $52 \pm 13$ | $48 \pm 11$ | $39 \pm 18$ | $55 \pm 21$ | $54 \pm 8$ |
| 10 | 0.5 | MPU-CE | 20 | $56 \pm 15$ | $55 \pm 16$ | $54 \pm 11$ | $44 \pm 12$ | $42 \pm 13$ | $31 \pm 18$ | $51 \pm 15$ | $48 \pm 8$ |

Table 10: Validation results in terms of F1-Score of XML-RoBERTa voters – Part I. Measures are micro-averaged among classes and then, for each setting, averaged among datasets and reported in percentage with their standard deviation. The MEAN column reports the performance averaged across settings with its standard deviation. The best results for each setting are reported in bold.

| E | Neg | Loss | m | 10 | 20 | 40 | 60 | 80 | 100 | Fully | MEAN |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.1 | MPN | 10 | $55 \pm 2$ | $74 \pm 14$ | $\mathbf{82 \pm 7}$ | $86 \pm 6$ | $91 \pm 2$ | $90 \pm 2$ | $93 \pm 3$ | $81 \pm 12$ |
| 5 | 0.1 | MPN | 20 | $55 \pm 2$ | $81 \pm 6$ | $77 \pm 11$ | $78 \pm 8$ | $80 \pm 11$ | $82 \pm 8$ | $89 \pm 5$ | $78 \pm 9$ |
| 5 | 0.1 | MPN-CE | 10 | $59 \pm 9$ | $75 \pm 1$ | $77 \pm 2$ | $82 \pm 2$ | $87 \pm 2$ | $87 \pm 3$ | $93 \pm 3$ | $80 \pm 10$ |
| 5 | 0.1 | MPN-CE | 20 | $64 \pm 9$ | $75 \pm 5$ | $75 \pm 7$ | $81 \pm 6$ | $83 \pm 4$ | $85 \pm 5$ | $92 \pm 4$ | $79 \pm 8$ |
| 5 | 0.1 | MPU | 10 | $67 \pm 12$ | $62 \pm 8$ | $64 \pm 10$ | $49 \pm 3$ | $44 \pm 1$ | $43 \pm 6$ | $36 \pm 7$ | $52 \pm 11$ |
| 5 | 0.1 | MPU | 20 | $67 \pm 15$ | $65 \pm 14$ | $59 \pm 11$ | $57 \pm 8$ | $45 \pm 8$ | $43 \pm 9$ | $37 \pm 7$ | $53 \pm 10$ |
| 5 | 0.1 | MPU-CE | 10 | $60 \pm 11$ | $62 \pm 14$ | $56 \pm 17$ | $51 \pm 9$ | $40 \pm 14$ | $41 \pm 2$ | $34 \pm 0$ | $49 \pm 9$ |
| 5 | 0.1 | MPU-CE | 20 | $56 \pm 13$ | $55 \pm 16$ | $53 \pm 14$ | $47 \pm 8$ | $41 \pm 6$ | $28 \pm 3$ | $32 \pm 1$ | $45 \pm 10$ |
| 5 | 0.3 | MPN | 10 | $59 \pm 2$ | $77 \pm 10$ | $\mathbf{82 \pm 8}$ | $86 \pm 5$ | $90 \pm 3$ | $90 \pm 2$ | $92 \pm 2$ | $\mathbf{82 \pm 10}$ |
| 5 | 0.3 | MPN | 20 | $59 \pm 5$ | $\mathbf{82 \pm 8}$ | $76 \pm 11$ | $79 \pm 9$ | $80 \pm 12$ | $82 \pm 8$ | $89 \pm 4$ | $78 \pm 8$ |
| 5 | 0.3 | MPN-CE | 10 | $59 \pm 9$ | $76 \pm 2$ | $78 \pm 1$ | $84 \pm 2$ | $86 \pm 1$ | $87 \pm 3$ | $92 \pm 3$ | $80 \pm 9$ |
| 5 | 0.3 | MPN-CE | 20 | $63 \pm 9$ | $76 \pm 6$ | $75 \pm 7$ | $81 \pm 5$ | $83 \pm 5$ | $85 \pm 5$ | $91 \pm 5$ | $79 \pm 8$ |
| 5 | 0.3 | MPU | 10 | $67 \pm 13$ | $54 \pm 2$ | $49 \pm 8$ | $45 \pm 5$ | $42 \pm 4$ | $41 \pm 9$ | $43 \pm 10$ | $49 \pm 8$ |
| 5 | 0.3 | MPU | 20 | $67 \pm 15$ | $64 \pm 16$ | $60 \pm 12$ | $48 \pm 3$ | $47 \pm 3$ | $43 \pm 10$ | $41 \pm 5$ | $53 \pm 9$ |
| 5 | 0.3 | MPU-CE | 10 | $61 \pm 12$ | $61 \pm 16$ | $57 \pm 16$ | $54 \pm 10$ | $45 \pm 16$ | $40 \pm 8$ | $34 \pm 4$ | $50 \pm 9$ |
| 5 | 0.3 | MPU-CE | 20 | $58 \pm 12$ | $56 \pm 16$ | $52 \pm 14$ | $46 \pm 9$ | $42 \pm 7$ | $26 \pm 1$ | $32 \pm 0$ | $44 \pm 10$ |
| 5 | 0.5 | MPN | 10 | $48 \pm 10$ | $76 \pm 11$ | $81 \pm 8$ | $86 \pm 5$ | $\mathbf{92 \pm 3}$ | $91 \pm 1$ | $92 \pm 2$ | $81 \pm 14$ |
| 5 | 0.5 | MPN | 20 | $57 \pm 1$ | $\mathbf{82 \pm 7}$ | $77 \pm 11$ | $79 \pm 9$ | $80 \pm 12$ | $82 \pm 8$ | $89 \pm 5$ | $78 \pm 9$ |
| 5 | 0.5 | MPN-CE | 10 | $60 \pm 8$ | $76 \pm 2$ | $76 \pm 2$ | $82 \pm 3$ | $86 \pm 2$ | $87 \pm 3$ | $93 \pm 4$ | $80 \pm 9$ |
| **5** | **0.5** | **MPN-CE** | **20** | $65 \pm 8$ | $76 \pm 6$ | $75 \pm 7$ | $81 \pm 5$ | $82 \pm 4$ | $85 \pm 4$ | $92 \pm 4$ | $80 \pm 7$ |
| 5 | 0.5 | MPU | 10 | $66 \pm 15$ | $54 \pm 1$ | $52 \pm 3$ | $48 \pm 4$ | $42 \pm 4$ | $39 \pm 7$ | $35 \pm 2$ | $48 \pm 9$ |
| 5 | 0.5 | MPU | 20 | $\mathbf{70 \pm 18}$ | $67 \pm 11$ | $61 \pm 12$ | $46 \pm 2$ | $46 \pm 4$ | $43 \pm 9$ | $41 \pm 6$ | $53 \pm 10$ |
| 5 | 0.5 | MPU-CE | 10 | $59 \pm 12$ | $60 \pm 13$ | $57 \pm 17$ | $53 \pm 11$ | $51 \pm 11$ | $38 \pm 11$ | $42 \pm 10$ | $51 \pm 7$ |
| 5 | 0.5 | MPU-CE | 20 | $59 \pm 13$ | $56 \pm 14$ | $53 \pm 13$ | $49 \pm 9$ | $41 \pm 9$ | $32 \pm 6$ | $46 \pm 12$ | $48 \pm 8$ |

Table 11: Validation results in terms of F1-Score of XML-RoBERTa voters – Part II. Measures are micro-averaged among classes and then, for each setting, averaged among datasets and reported in percentage with their standard deviation. The MEAN column reports the performance averaged across settings with its standard deviation. The best results for each setting are reported in bold. Selected hyperparameters are in bold.

| E | Loss | 10 | 20 | 40 | 60 | 80 | 100 | Fully | MEAN |
|---|---|---|---|---|---|---|---|---|---|
| 15 | Conf-MPU-CE | $77 \pm 1$ | $80 \pm 6$ | $82 \pm 3$ | $81 \pm 4$ | $79 \pm 7$ | $79 \pm 9$ | $90 \pm 1$ | $81 \pm 4$ |
| 15 | Conf-MPU | $\mathbf{81 \pm 1}$ | $83 \pm 5$ | $81 \pm 3$ | $83 \pm 4$ | $80 \pm 8$ | $78 \pm 8$ | $89 \pm 3$ | $82 \pm 3$ |
| 15 | MPN-CE | $72 \pm 5$ | $77 \pm 3$ | $80 \pm 2$ | $82 \pm 4$ | $82 \pm 7$ | $80 \pm 8$ | $92 \pm 1$ | $81 \pm 5$ |
| 15 | MPN | $78 \pm 2$ | $81 \pm 5$ | $82 \pm 2$ | $82 \pm 3$ | $82 \pm 7$ | $80 \pm 7$ | $90 \pm 3$ | $82 \pm 3$ |
| 15 | MPU-CE | $64 \pm 12$ | $66 \pm 12$ | $59 \pm 14$ | $51 \pm 14$ | $44 \pm 14$ | $29 \pm 18$ | $43 \pm 21$ | $51 \pm 12$ |
| 15 | MPU | $59 \pm 26$ | $66 \pm 23$ | $62 \pm 23$ | $60 \pm 20$ | $58 \pm 17$ | $47 \pm 20$ | $58 \pm 33$ | $59 \pm 5$ |
| 20 | Conf-MPU-CE | $77 \pm 2$ | $78 \pm 8$ | $81 \pm 2$ | $81 \pm 4$ | $81 \pm 7$ | $80 \pm 9$ | $91 \pm 1$ | $81 \pm 4$ |
| 20 | Conf-MPU | $77 \pm 2$ | $82 \pm 6$ | $83 \pm 2$ | $82 \pm 5$ | $80 \pm 6$ | $79 \pm 8$ | $90 \pm 3$ | $82 \pm 4$ |
| 20 | MPN-CE | $69 \pm 5$ | $74 \pm 6$ | $80 \pm 1$ | $83 \pm 5$ | $82 \pm 6$ | $\mathbf{83 \pm 9}$ | $\mathbf{94 \pm 1}$ | $80 \pm 7$ |
| **20** | **MPN** | $80 \pm 0$ | $83 \pm 6$ | $\mathbf{85 \pm 3}$ | $\mathbf{84 \pm 4}$ | $82 \pm 6$ | $82 \pm 7$ | $90 \pm 3$ | $\mathbf{84 \pm 3}$ |
| 20 | MPU-CE | $63 \pm 11$ | $63 \pm 17$ | $58 \pm 15$ | $51 \pm 12$ | $43 \pm 15$ | $27 \pm 18$ | $38 \pm 19$ | $49 \pm 12$ |
| 20 | MPU | $65 \pm 18$ | $64 \pm 26$ | $74 \pm 9$ | $68 \pm 10$ | $63 \pm 11$ | $60 \pm 12$ | $61 \pm 29$ | $65 \pm 4$ |
| 25 | Conf-MPU-CE | $76 \pm 3$ | $81 \pm 5$ | $82 \pm 3$ | $82 \pm 4$ | $82 \pm 7$ | $81 \pm 9$ | $91 \pm 1$ | $82 \pm 4$ |
| 25 | Conf-MPU | $79 \pm 2$ | $82 \pm 6$ | $84 \pm 3$ | $82 \pm 5$ | $79 \pm 7$ | $80 \pm 8$ | $90 \pm 3$ | $82 \pm 3$ |
| 25 | MPN-CE | $70 \pm 5$ | $75 \pm 4$ | $78 \pm 1$ | $82 \pm 4$ | $\mathbf{83 \pm 6}$ | $\mathbf{83 \pm 9}$ | $93 \pm 1$ | $81 \pm 6$ |
| 25 | MPN | $77 \pm 3$ | $\mathbf{84 \pm 6}$ | $84 \pm 2$ | $\mathbf{84 \pm 3}$ | $83 \pm 7$ | $82 \pm 8$ | $91 \pm 2$ | $\mathbf{84 \pm 3}$ |
| 25 | MPU-CE | $65 \pm 12$ | $61 \pm 17$ | $58 \pm 16$ | $50 \pm 15$ | $44 \pm 16$ | $27 \pm 18$ | $36 \pm 20$ | $49 \pm 12$ |
| 25 | MPU | $70 \pm 13$ | $67 \pm 23$ | $76 \pm 6$ | $68 \pm 8$ | $65 \pm 9$ | $61 \pm 9$ | $63 \pm 26$ | $67 \pm 4$ |

Table 12: Validation results in terms of F1-Score of CuPUL curriculum train with Stacked Embeddings. Measures are micro-averaged among classes and then, for each setting, averaged among datasets and reported in percentage with their standard deviation. The MEAN column reports the performance averaged across settings with its standard deviation. The best results for each setting are reported in bold. Selected hyperparameters are in bold.

| E | sE | Loss | LR | 10 | 20 | 40 | 60 | 80 | 100 | Fully | MEAN |
|---|----|------|----|----|----|----|----|----|-----|-------|------|
| 1 | 1 | Conf-MPU | 1 | 28 ± 24 | 36 ± 31 | 40 ± 27 | 42 ± 33 | 43 ± 37 | 42 ± 38 | 44 ± 40 | 39 ± 5 |
| 1 | 1 | Conf-MPU | 3 | 51 ± 29 | 50 ± 33 | 57 ± 28 | 57 ± 32 | 58 ± 30 | 53 ± 35 | 56 ± 36 | 55 ± 3 |
| 1 | 1 | Conf-MPU-CE | 1 | 29 ± 25 | 43 ± 33 | 39 ± 28 | 36 ± 34 | 38 ± 35 | 39 ± 34 | 44 ± 41 | 38 ± 4 |
| 1 | 1 | Conf-MPU-CE | 3 | 67 ± 15 | 68 ± 13 | 61 ± 17 | 60 ± 18 | 60 ± 17 | 66 ± 14 | 74 ± 17 | 65 ± 5 |
| 1 | 1 | MPU-CE | 1 | 27 ± 23 | 37 ± 31 | 38 ± 26 | 31 ± 28 | 30 ± 27 | 26 ± 23 | 33 ± 30 | 32 ± 4 |
| 1 | 1 | MPU-CE | 3 | 59 ± 15 | 57 ± 6 | 46 ± 7 | 53 ± 10 | 37 ± 5 | 34 ± 6 | 27 ± 12 | 45 ± 11 |
| 1 | 2 | MPU-CE | 1 | 47 ± 17 | 53 ± 11 | 48 ± 16 | 46 ± 13 | 45 ± 10 | 43 ± 5 | 35 ± 5 | 45 ± 5 |
| 1 | 2 | MPU-CE | 3 | 53 ± 12 | 61 ± 6 | 54 ± 9 | 52 ± 11 | 35 ± 7 | 42 ± 7 | 30 ± 10 | 47 ± 10 |
| 5 | 1 | Conf-MPU | 1 | 70 ± 8 | 76 ± 9 | 69 ± 14 | 73 ± 14 | 74 ± 17 | 79 ± 12 | 79 ± 14 | 74 ± 3 |
| 5 | 1 | Conf-MPU | 3 | 69 ± 0 | 80 ± 4 | 79 ± 7 | 86 ± 3 | 83 ± 7 | 89 ± 1 | 92 ± 3 | 82 ± 7 |
| 5 | 1 | Conf-MPU-CE | 1 | 69 ± 5 | 78 ± 10 | 74 ± 11 | 76 ± 7 | 76 ± 8 | 80 ± 5 | 91 ± 3 | 78 ± 6 |
| 5 | 1 | Conf-MPU-CE | 3 | 70 ± 7 | 80 ± 2 | 81 ± 1 | 84 ± 3 | 88 ± 1 | 87 ± 1 | 93 ± 2 | 83 ± 6 |
| 5 | 1 | MPU-CE | 1 | 56 ± 9 | 55 ± 11 | 52 ± 14 | 49 ± 13 | 46 ± 12 | 39 ± 7 | 39 ± 13 | 48 ± 6 |
| 5 | 1 | MPU-CE | 3 | 52 ± 12 | 67 ± 3 | 57 ± 8 | 57 ± 1 | 47 ± 4 | 50 ± 5 | 42 ± 11 | 53 ± 7 |
| 5 | 2 | Conf-MPU | 1 | 66 ± 9 | 79 ± 4 | 81 ± 3 | 88 ± 1 | 85 ± 7 | 90 ± 2 | 92 ± 3 | 83 ± 8 |
| 5 | 2 | Conf-MPU | 3 | 61 ± 3 | **84 ± 0** | 80 ± 6 | 87 ± 3 | 89 ± 1 | **91 ± 1** | 94 ± 2 | **84 ± 9** |
| 5 | 2 | Conf-MPU-CE | 1 | **73 ± 14** | 81 ± 3 | 81 ± 4 | 78 ± 8 | 83 ± 5 | 88 ± 0 | 93 ± 3 | 83 ± 6 |
| 5 | 2 | Conf-MPU-CE | 3 | 64 ± 14 | 77 ± 3 | 81 ± 4 | 81 ± 7 | 84 ± 5 | 90 ± 0 | 94 ± 3 | 82 ± 8 |
| 5 | 2 | MPU-CE | 1 | 65 ± 7 | 63 ± 3 | 56 ± 4 | 53 ± 5 | 49 ± 10 | 44 ± 6 | 57 ± 12 | 55 ± 7 |
| 5 | 2 | MPU-CE | 3 | 63 ± 3 | 60 ± 13 | 50 ± 17 | 66 ± 5 | 59 ± 10 | 44 ± 21 | 65 ± 21 | 58 ± 7 |
| 10 | 1 | Conf-MPU | 1 | 70 ± 5 | 79 ± 5 | 80 ± 5 | 85 ± 4 | 84 ± 6 | 88 ± 3 | 93 ± 2 | 83 ± 6 |
| 10 | 1 | Conf-MPU | 3 | 65 ± 10 | 78 ± 2 | 82 ± 4 | 89 ± 0 | 87 ± 4 | 89 ± 2 | 92 ± 3 | 83 ± 8 |
| **10** | **1** | **Conf-MPU-CE** | **1** | **73 ± 11** | 82 ± 2 | 81 ± 4 | 80 ± 8 | 85 ± 4 | 89 ± 0 | 92 ± 5 | 83 ± 5 |
| 10 | 1 | Conf-MPU-CE | 3 | 62 ± 16 | 75 ± 6 | **86 ± 0** | 83 ± 5 | 86 ± 2 | 90 ± 0 | 94 ± 3 | 82 ± 9 |
| 10 | 1 | MPU-CE | 1 | 62 ± 8 | 63 ± 11 | 54 ± 8 | 51 ± 10 | 51 ± 6 | 47 ± 3 | 60 ± 13 | 55 ± 5 |
| 10 | 1 | MPU-CE | 3 | 57 ± 7 | 63 ± 17 | 52 ± 20 | 51 ± 20 | 60 ± 1 | 51 ± 11 | 56 ± 31 | 56 ± 4 |
| 10 | 2 | Conf-MPU | 1 | 61 ± 8 | 83 ± 2 | 85 ± 1 | 86 ± 3 | 91 ± 1 | 90 ± 1 | 92 ± 4 | **84 ± 9** |
| 10 | 2 | Conf-MPU | 3 | 54 ± 9 | 83 ± 2 | **86 ± 0** | **90 ± 0** | **92 ± 1** | **91 ± 1** | 91 ± 4 | **84 ± 12** |
| 10 | 2 | Conf-MPU-CE | 1 | 67 ± 17 | 77 ± 3 | 78 ± 7 | 79 ± 8 | 85 ± 4 | **91 ± 1** | **95 ± 2** | 82 ± 8 |
| 10 | 2 | Conf-MPU-CE | 3 | 57 ± 17 | 80 ± 3 | 82 ± 1 | 80 ± 8 | 87 ± 1 | 90 ± 0 | **95 ± 2** | 82 ± 11 |
| 10 | 2 | MPU-CE | 1 | 60 ± 13 | 55 ± 9 | 62 ± 3 | 54 ± 11 | 55 ± 9 | 46 ± 6 | 67 ± 19 | 57 ± 6 |
| 10 | 2 | MPU-CE | 3 | 67 ± 8 | 62 ± 14 | 67 ± 7 | 63 ± 15 | 72 ± 5 | 59 ± 17 | 63 ± 14 | 65 ± 4 |

Table 13: Validation results in terms of F1-Score of CuPUL curriculum train with XLM-RoBERTa. Measures are micro-averaged among classes and then, for each setting, averaged among datasets and reported in percentage with their standard deviation. The MEAN column reports the performance averaged across settings with its standard deviation. The best results for each setting are reported in bold. Selected hyperparameters are in bold. Learning rate is in the scale of $1e-5$