## A Appendices

### A.1 Hyperparameters

We tuned the batch size and the learning rate using development sets in four languages,[14] and then fixed these hyperparameters for all other languages in each model. The batch size was 1 sentence in low-resource scenarios (in baseline `LSup` and fine-tuning of `RaRe`), and to 100 sentences, in high-resource settings (`HSup` and the pretraining phase of `RaRe`). The learning rate was set to 0.001 and 0.01 for the high-resource and low-resource baseline models, respectively, and to 0.005, 0.0005 for the pretraining and fine-tuning phase of `RaRe` based on development results for the four languages. For CoNLL datasets, we had to decrease the batch size of the pre-training phase from 100 to 20 (because of GPU memory issues).

### A.2 Cross-lingual Word Embeddings

We experimented with Wiki and CommonCrawl monolingual embeddings from `fastText` (Bojanowski et al., 2017). Each of the 41 languages is mapped to English embedding space using three methods from `MUSE`: 1) supervised with bilingual dictionaries; 2) seeding using identical character sequences; and 3) unsupervised training using adversarial learning (Lample et al., 2018). The cross-lingual mappings are evaluated by precision at $k = 1$. The resulting cross-lingual embeddings are then used in NER direct transfer in a leave-one-out setting for the 41 languages ($41 \times 40$ transfers), and we report the mean $F_1$ in Table 3. CommonCrawl doesn't perform well in bilingual induction despite having larger text corpora, and underperforms in direct transfer NER. It is also evident that using identical character strings instead of a bilingual dictionary as the seed for learning a supervised bilingual mapping barely effects the performance. This finding also applies to few-shot learning over larger ensembles: running `RaRe` over 40 source languages achieves an average $F_1$ of 77.9 when using embeddings trained with a dictionary, versus 76.9 using string identity instead. For this reason we have used the string identity method in the paper (e.g., Table 4), providing greater portability to language pairs without a bilingual dictionary. Experiments with unsupervised mappings performed substantially worse than supervised methods, and so we didn't explore these further.

---

[14] Afrikaans, Arabic, Bulgarian and Bengali.

| | Unsup | | IdentChr | | Sup | |
|---|---|---|---|---|---|---|
| | crawl | wiki | crawl | wiki | crawl | wiki |
| word translation accuracy | | | | | | |
| | 34 | 24 | 43 | **53** | 50 | **54** |
| average $F_1$ in direct transfer | | | | | | |
| | 26 | 21 | 37 | **44** | 39 | **45** |

**Table 3:** The effect of the choice of monolingual word embeddings (Common Crawl and Wikipedia), and their cross-lingual mapping on NER direct transfer.

### A.3 Direct Transfer Results

In Figure 5 the performance of an NER model trained in a high-resource setting on a source language applied on the other 40 target languages (leave-one-out) is shown. An interesting finding is that symmetry does not always hold (e.g. id vs. ms or fa vs. ar).

### A.4 Detailed Low-resource Results

The result of applying baselines, proposed models and their variations, and unsupervised transfer model of Xie et al. (2018) are shown in Table 4.

**Figure 5:** The direct transfer performance of a source NER model trained in a high-resource setting applied on the other 40 target languages, and evaluated in terms of phrase-level $F_1$.

| | #train (k) | #test (k) | BiDic.P@1 | Supervised | | | | | | Unsupervised | | | | | | Oracle |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | HSup | LSup | RaRe t1 | RaRe t10 | RaRe all | $\text{BEA}_{sup}^{ent}$ t10 | $\text{RaRe}_{uns}$ | BWET | $\text{BEA}_{uns\times2}^{ent}$ t10 | $\text{BEA}_{uns}^{ent}$ | $\text{BEA}_{uns}^{tok}$ | $\text{MV}^{tok}$ | |
| af | 5 | 1 | 36 | 84 | 59 | 73 | 79 | 79 | 80 | 76 | 64 | 79 | 79 | 74 | 75 | 80 |
| ar | 20 | 10 | 46 | 88 | 64 | 71 | 74 | 74 | 65 | 26 | 19 | 54 | 45 | 54 | 12 | 56 |
| bg | 20 | 10 | 55 | 90 | 61 | 80 | 81 | 81 | 81 | 5 | 51 | 81 | 65 | 54 | 4 | 76 |
| bn | 10 | 1 | 1 | 95 | 70 | 68 | 74 | 74 | 69 | 65 | 36 | 67 | 66 | 60 | 56 | 63 |
| bs | 15 | 1 | 30 | 92 | 63 | 80 | 79 | 80 | 78 | 76 | 52 | 80 | 78 | 77 | 69 | 82 |
| ca | 20 | 10 | 70 | 91 | 62 | 82 | 86 | 84 | 86 | 80 | 62 | 85 | 80 | 79 | 72 | 83 |
| cs | 20 | 10 | 64 | 90 | 62 | 77 | 78 | 75 | 78 | 73 | 59 | 77 | 75 | 72 | 71 | 78 |
| da | 20 | 10 | 68 | 90 | 62 | 77 | 81 | 81 | 82 | 79 | 68 | 83 | 82 | 79 | 78 | 80 |
| de | 20 | 10 | 73 | 86 | 58 | 73 | 74 | 73 | 72 | 69 | 63 | 72 | 71 | 64 | 68 | 70 |
| el | 20 | 10 | 55 | 89 | 61 | 67 | 67 | 67 | 54 | 13 | 45 | 49 | 43 | 34 | 13 | 45 |
| en | 20 | 10 | — | 81 | 47 | 64 | 65 | 64 | 65 | 58 | — | 63 | 61 | 57 | 56 | 61 |
| es | 20 | 10 | 83 | 90 | 63 | 83 | 84 | 84 | 85 | 76 | 62 | 85 | 81 | 76 | 73 | 84 |
| et | 15 | 10 | 41 | 90 | 64 | 73 | 77 | 77 | 78 | 72 | 58 | 78 | 78 | 71 | 73 | 75 |
| fa | 20 | 10 | 33 | 93 | 74 | 78 | 81 | 79 | 69 | 30 | 16 | 65 | 50 | 52 | 15 | 60 |
| fi | 20 | 10 | 58 | 89 | 67 | 78 | 80 | 80 | 81 | 76 | 68 | 81 | 80 | 69 | 77 | 78 |
| fr | 20 | 10 | 82 | 88 | 57 | 81 | 81 | 80 | 84 | 75 | 59 | 83 | 79 | 73 | 71 | 80 |
| he | 20 | 10 | 52 | 85 | 53 | 61 | 61 | 60 | 55 | 40 | 26 | 54 | 54 | 46 | 34 | 50 |
| hi | 5 | 1 | 29 | 85 | 68 | 64 | 74 | 73 | 68 | 48 | 27 | 64 | 61 | 58 | 35 | 54 |
| hr | 20 | 10 | 48 | 89 | 61 | 74 | 79 | 78 | 80 | 76 | 49 | 80 | 79 | 77 | 73 | 78 |
| hu | 20 | 10 | 64 | 90 | 59 | 75 | 79 | 78 | 80 | 71 | 55 | 79 | 79 | 69 | 73 | 76 |
| id | 20 | 10 | 68 | 91 | 67 | 82 | 83 | 81 | 75 | 59 | 62 | 73 | 67 | 61 | 62 | 79 |
| it | 20 | 10 | 77 | 89 | 60 | 80 | 81 | 80 | 82 | 75 | 59 | 81 | 78 | 76 | 72 | 79 |
| lt | 10 | 10 | 26 | 86 | 62 | 72 | 79 | 80 | 79 | 76 | 48 | 80 | 80 | 75 | 77 | 74 |
| lv | 10 | 10 | 31 | 91 | 68 | 70 | 75 | 75 | 69 | 68 | 40 | 69 | 69 | 67 | 65 | 66 |
| mk | 10 | 1 | 50 | 91 | 67 | 79 | 82 | 81 | 80 | 4 | 38 | 79 | 66 | 48 | 3 | 75 |
| ms | 20 | 1 | 48 | 91 | 66 | 78 | 80 | 78 | 74 | 69 | 62 | 68 | 67 | 63 | 68 | 74 |
| nl | 20 | 10 | 76 | 89 | 59 | 78 | 80 | 80 | 81 | 77 | 63 | 82 | 81 | 78 | 76 | 79 |
| no | 20 | 10 | 67 | 90 | 65 | 79 | 82 | 81 | 83 | 79 | 59 | 83 | 83 | 77 | 79 | 79 |
| pl | 20 | 10 | 66 | 89 | 61 | 76 | 79 | 78 | 81 | 73 | 63 | 82 | 80 | 77 | 76 | 78 |
| pt | 20 | 10 | 80 | 90 | 59 | 79 | 81 | 80 | 82 | 77 | 65 | 82 | 77 | 74 | 70 | 82 |
| ro | 20 | 10 | 67 | 92 | 66 | 80 | 82 | 82 | 80 | 76 | 46 | 78 | 76 | 74 | 67 | 77 |
| ru | 20 | 10 | 59 | 86 | 53 | 73 | 71 | 71 | 56 | 10 | 38 | 53 | 40 | 36 | 11 | 61 |
| sk | 20 | 10 | 52 | 91 | 62 | 76 | 79 | 79 | 80 | 74 | 50 | 79 | 76 | 76 | 71 | 79 |
| sl | 15 | 10 | 47 | 92 | 64 | 76 | 80 | 80 | 79 | 76 | 58 | 79 | 78 | 76 | 73 | 78 |
| sq | 5 | 1 | 37 | 88 | 69 | 79 | 84 | 84 | 83 | 82 | 59 | 83 | 84 | 76 | 79 | 79 |
| sv | 20 | 10 | 61 | 93 | 69 | 83 | 83 | 84 | 82 | 77 | 60 | 79 | 80 | 69 | 76 | 84 |
| ta | 15 | 1 | 7 | 84 | 54 | 44 | 53 | 53 | 46 | 35 | 12 | 39 | 42 | 25 | 29 | 38 |
| tl | 10 | 1 | 20 | 93 | 66 | 75 | 82 | 80 | 78 | 65 | 60 | 62 | 60 | 57 | 52 | 76 |
| tr | 20 | 10 | 61 | 90 | 61 | 75 | 77 | 77 | 77 | 70 | 53 | 77 | 76 | 67 | 67 | 71 |
| uk | 20 | 10 | 45 | 89 | 60 | 70 | 78 | 79 | 70 | 5 | 35 | 64 | 58 | 49 | 6 | 60 |
| vi | 20 | 10 | 54 | 88 | 55 | 64 | 72 | 72 | 61 | 58 | 53 | 56 | 55 | 48 | 47 | 56 |
| $\mu$ | — | — | — | 89.2 | 62.1 | 74.3 | 77.4 | 76.9 | 74.8 | 60.2 | 50.5 | 72.8 | 69.7 | 64.5 | 56.7 | 71.6 |
| $\sigma$ | — | — | — | 2.8 | 5.2 | 7.3 | 6.4 | 6.4 | 9.6 | 24.1 | 14.7 | 11.5 | 12.6 | 13.7 | 25 | 11.5 |

**Table 4:** The size of training and test sets (development set size equals test set size) in thousand sentences, and the precision at 1 for Bilingual dictionaries induced from mapping languages to the English embedding space (using identical characters) is shown (`BiDic.P@1`). $F_1$ scores on the test set, comparing baseline supervised models (`HSup`, `LSup`), multilingual transfer from top $k$ source languages (`RaRe`, 5 runs, $k = 1, 10, 40$), an unsupervised `RaRe` with uniform expertise and no fine-tuning (`RaRe`$_{uns}$), and aggregation methods: majority voting (`MV`$^{tok}$), `BEA`$_{uns}^{tok}$ and `BEA`$_{uns}^{ent}$ (Bayesian aggregation in token- and entity-level), and the oracle single best annotation (`Oracle`). We also compare with `BWET` (Xie et al., 2018), an unsupervised transfer model with state-of-the-art on CoNLL NER datasets. The mean and standard deviation over all 41 languages, $\mu, \sigma$, are also reported.