

| EMEA             | AUC          | MacF         | Science          | AUC          | MacF         | Subs             | AUC          | MacF         |
|------------------|--------------|--------------|------------------|--------------|--------------|------------------|--------------|--------------|
| ALLFEATURES      | 79.60        | 71.64        | ALLFEATURES      | 73.91        | 47.54        | ALLFEATURES      | 69.26        | 52.78        |
| –Token:L-PSDBin  | 77.09        | 70.50        | –Token:L-PSDBin  | 76.26        | 53.69        | –Type:NgramProb  | 69.13        | 53.33        |
| –Type:RelFreq    | 78.43        | 72.19        | –Token:G-PSD     | 77.04        | 53.56        | –Token:G-PSDBin  | 70.23        | 54.72        |
| –Token:G-PSD     | <b>79.66</b> | 72.11        | –Token:G-PSDBin  | 77.44        | 54.54        | –Token:CtxCnt    | 71.23        | 58.35        |
| –Type:Context    | <b>79.66</b> | 72.45        | –Token:L-PSD     | 77.85        | 56.05        | –Token:L-PSDBin  | 72.07        | 57.85        |
| –Token:Ctx%      | 78.91        | <b>73.37</b> | –Token:PSDRatio  | 77.92        | <b>57.34</b> | –Token:G-PSD     | 72.17        | 57.33        |
| –Type:TopicSim   | 78.05        | 71.33        | –Token:CtxCnt    | 77.85        | 54.42        | –Type:TopicSim   | <b>72.31</b> | 58.41        |
| –Token:CtxCnt    | 76.90        | 71.72        | –Type:Context    | <b>78.17</b> | 55.45        | –Token:Ctx%      | 72.17        | 56.17        |
| –Token:L-PSD     | 76.03        | 73.35        | –Token:Ctx%      | 78.06        | 55.04        | –Token:NgramProb | 71.35        | <b>59.26</b> |
| –Type:NgramProb  | 73.32        | 69.54        | –Type:TopicSim   | 77.83        | 54.57        | –Token:PSDRatio  | 70.33        | 46.88        |
| –Token:G-PSDBin  | 74.41        | 69.76        | –Token:NgramProb | 76.98        | 51.02        | –Token:L-PSD     | 69.05        | 53.31        |
| –Token:NgramProb | 69.78        | 68.89        | –Type:RelFreq    | 74.25        | 49.57        | –Type:RelFreq    | 65.25        | 48.22        |
| –Token:PSDRatio  | 48.38        | 3.45         | –Type:NgramProb  | 50.00        | 0.00         | –Type:Context    | 50.00        | 0.00         |

Table 5: Feature ablation results for all three corpora. Selection criteria is AUC, but Macro-F is presented for completeness. Feature selection is run independently on each of the three datasets. The features toward the *bottom* were the first selected.

|                | AUC                     | Macro-F                 | Micro-F                 |
|----------------|-------------------------|-------------------------|-------------------------|
| <b>EMEA</b>    |                         |                         |                         |
| TYPEONLY       | 71.43 $\pm$ 0.94        | 52.62 $\pm$ 3.41        | 38.67 $\pm$ 1.35        |
| TOKENONLY      | <b>73.75</b> $\pm$ 1.11 | <b>67.77</b> $\pm$ 4.18 | 45.49 $\pm$ 3.96        |
| ALLFEATURES    | 72.19 $\pm$ 4.07        | 67.26 $\pm$ 7.88        | <b>49.29</b> $\pm$ 3.55 |
| XV-ALLFEATURES | 79.60 $\pm$ 1.20        | 71.64 $\pm$ 1.83        | 46.83 $\pm$ 0.62        |
| <b>Science</b> |                         |                         |                         |
| TYPEONLY       | <b>75.19</b> $\pm$ 0.89 | <b>51.53</b> $\pm$ 2.55 | 37.14 $\pm$ 4.41        |
| TOKENONLY      | 71.24 $\pm$ 1.45        | 47.27 $\pm$ 1.11        | 40.48 $\pm$ 1.84        |
| ALLFEATURES    | 74.14 $\pm$ 0.93        | 48.86 $\pm$ 3.94        | <b>43.20</b> $\pm$ 3.16 |
| XV-ALLFEATURES | 73.91 $\pm$ 0.66        | 47.54 $\pm$ 1.52        | 40.22 $\pm$ 1.03        |
| <b>Subs</b>    |                         |                         |                         |
| TYPEONLY       | 60.90 $\pm$ 1.47        | 39.21 $\pm$ 14.78       | 24.77 $\pm$ 2.78        |
| TOKENONLY      | <b>62.00</b> $\pm$ 1.16 | 49.74 $\pm$ 6.30        | <b>42.95</b> $\pm$ 3.92 |
| ALLFEATURES    | 60.12 $\pm$ 2.11        | <b>50.16</b> $\pm$ 8.63 | 38.56 $\pm$ 5.20        |
| XV-ALLFEATURES | 69.26 $\pm$ 0.60        | 52.78 $\pm$ 1.96        | 43.85 $\pm$ 0.90        |

Table 6: Cross-domain test results on the SENSESPOTTING task. Two standard deviations are shown in small type. Only AUC, Macro-F and Micro-F are shown for brevity.

AUC as the measure on which to ablate. It’s quite clear that for Science, all the useful information is in the type-level features, a result that echoes what we saw in the previous section. While for EMEA and Subs, both type- and token-level features play a significant role. Considering the six most useful features in each domain, the ones that pop out as frequently most useful are the global PSD features, the ngram probability features (either type- or token-based), the relative frequency features and the context features.

## 6.5 Cross-Domain Training

One disadvantage to the previous method for evaluating the SENSESPOTTING task is that it requires parallel data in a new domain. Suppose we have *no* parallel data in the new domain at all, yet still want to attack the SENSESPOTTING task. One option is

to train a system on domains for which we *do* have parallel data, and then apply it in a new domain. This is precisely the setting we explore in this section. Now, instead of performing cross-validation in a single domain (for instance, Science), we take the union of *all* of the training data in the other domains (e.g., EMEA and Subs), train a classifier, and then apply it to Science. This classifier will almost certainly be worse than one trained on NEW (Science) but does not require *any* parallel data in that domain. (Hyperparameters are chosen by development data from the OLD union.)

The results of this experiment are shown in Table 6. We include results for TOKENONLY, TYPEONLY and ALLFEATURES; all of these are trained in the cross-domain setting. To ease comparison to the results that do not suffer from domain shift, we also present “XV-ALLFEATURES”, which are results copied from Table 4 in which parallel data from NEW is used. Overall, there is a drop of about 7.3% absolute in AUC, moving from XV-ALLFEATURES to ALLFEATURES, including a small improvement in Science (likely because Science is markedly smaller than Subs, and “more difficult” than EMEA with many word types).

## 6.6 Detecting Most Frequent Sense Changes

We define a second, related task: MOSTFRE-QSENSECHANGE. In this task, instead of predicting if a given word token has a sense which is brand new with respect to the old domain, we predict whether it is being used with a sense which is not the one that was observed *most frequently* in the old domain. In our EMEA, Science, and Subtitles data, 68.2%, 48.3%, and 69.6% of word tokens’ predominant sense changes.