

Politicians Are Also People: Mapping is All You Need

Clustering Entity Types in Cross-Domain Relation Classification Setups

David Peter Süle and Mie Jonasson and Nicklas Koch Rasmussen

BSc. in Data Science, IT University of Copenhagen

{dasy,miejo,nicra}@itu.dk

Abstract

Relation Extraction is an evolving field within natural language processing. As its last step, Relation Classification (RC) aims to identify the relation type to which two semantically related named entities belong. Cross-domain setups are especially challenging, even more so when domain-specific entity types are used. Research is scarce in the area and mostly focuses on using generic entity types or simply fine-tune the model on a single target domain. This might still offer challenges when annotated data is not accessible for fine-tuning.

In this paper we explore ways of clustering domain-specific named entity types to reduce cross-domain complexity and improve performance on previously unseen domains. We propose five different methods of grouping entity types and evaluate them in multi-domain and out-of-domain scenarios using our two new benchmarks. In conclusion, we find that all our entity mapping methods outperform the baseline in the out-of-domain setting, with the best performing model improving on the baseline by 8.6 percentage points in weighted F1.¹

1 Introduction

Relation extraction (RE) is an evolving field within natural language processing that aims to discover and classify relations between named entities in a given sentence. Despite out-of-domain generalization being an extremely desirable quality of RE models, research in this direction is lacking (Bassignana and Plank, 2022b). One complicating factor is that named entity types vary across domains and the most common solution is to fine-tune the model on target domain data. In many cases only limited amount of training data is available which poses another challenge. Research combining these two scenarios – *few-shot learning* and *domain shift* – is limited (Gao et al., 2019).

¹Our code can be found at: <https://github.com/itu.dk/dasy/2yp-project>

We suggest and evaluate methods for cross-domain clustering of named entity (NE) types to help relation classification (RC) models generalize better while preserving the necessary level of granularity. We work on the research question: *What is the performance impact of clustering domain-specific named entity types in cross-domain relation-classification setups and what benchmark can be established for future research?*

Our contributions are twofold:

- We propose two benchmarking methods for cross-domain RC models based on the publicly available CrossRE dataset (Bassignana and Plank, 2022a), matching two real life use cases, namely, multi-domain and out-of-domain evaluation.
- We evaluate five different NE type clustering methods with our proposed benchmarks.

2 Related Work

The CrossNER project (Liu et al., 2021) produced a human-annotated NER dataset of six diverse domains with domain-specific entity types. They highlight the lack of domain-specific entities in cross-domain benchmark datasets, leading to less generalizable evaluation results.

The CrossRE project (Bassignana and Plank, 2022a) builds upon the CrossNER data to bring a new cross-domain dataset for relation classification with an established baseline to be improved upon. The project notes the difficulty of the different label distributions across domains, and further highlights the need for research in the underexplored cross-domain dimension of RE.

3 Data

To test our proposed named entity grouping methods and the resulting relation classification model, we use the publicly available cross-domain

CrossRE dataset (Bassignana and Plank, 2022a) which itself is based upon the CrossNER project data (Liu et al., 2021). CrossNER consists of a curated set of sentences in six domains, namely AI, literature, music, natural science, news, and politics, each with domain-specific named entity labels and boundary annotations.

The two sources are Reuters News from the CoNLL-2003 shared task by Tjong Kim Sang and De Meulder (2003) (news domain), and Wikipedia. For details on the data collection, annotation processes, the vocabulary and data statistics we refer to Liu et al. (2021) and the published data.²

CrossRE uses the named entities identified in CrossNER (with minor changes) to hand-annotate semantic relations between them, using 17 domain-agnostic relation labels. The dataset models a low-resource scenario for domain-specific fine-tuning, with ~100 / ~350 / ~400 examples in the train/dev/test splits, respectively. This matches the original, resulting in over five thousand sentences.

Both entity and relation type distributions vary across domains; for statistics and annotation details we refer to Bassignana and Plank (2022a) and the published data.³ For the named entity and relation labels see Appendix A.⁴

4 Entity Type Clustering Methods

Examples and technical details for all clustering methods can be found in Appendices B through G.

4.1 Manual Grouping

For the first method, we started with a list of all entity types and attempted to create meaningful categories for subsets of the entities. We based our initial clusters on a grouping obtained from Elisa Bassignana (Appendix B.1). Throughout the process, we made individual categorizations, used the definitions of the entity types by Liu et al. (2021) to settle doubts, and discussed ambiguous cases.

4.2 Embedding Space Based Grouping

In our second method, we utilized word embedding vectors for algorithmic clustering using the google-news-300 (Google, 2013) word2vec model (Mikolov et al., 2013). This was a suitable choice as Liu et al. (2021) used the news as

²<https://github.com/zliucr/CrossNER>

³<https://github.com/mainlp/CrossRE>

⁴As we didn't modify the data set, we did not prepare a data statement (Bender and Friedman, 2018) and refer the reader to Appendix A of Bassignana and Plank (2022a) instead.

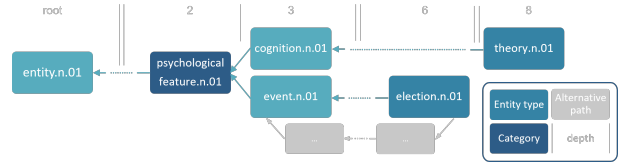


Figure 1: DAG structure of WordNet nouns. Both *theory* and *election* have *psychological feature* as their depth 2 hypernym, which is also their LCS.

the source domain. To obtain the embeddings, we simplified named entity labels consisting of two words (e.g. musical artist → musician) or replaced them with a specific example (e.g. programming language → javascript); see Appendix G for the full list of substitutions.

We applied a two-phase dimensionality reduction step consisting of PCA and UMAP (Uniform Manifold Approximation and Projection, McInnes et al., 2020). Using UMAP was motivated by our aim to more accurately capture fine detail and favour local information over preserving the global structure. To avoid the risk of picking up patterns from random noise we applied PCA before UMAP (Lee and Verleysen, 2007 and van der Maaten and Hinton, 2008). Finally, we used affinity propagation (Frey and Dueck, 2007) to cluster the results. This has the advantage of automatically determining the number of clusters, and being relatively stable across random seeds.

4.3 Topological Grouping

We use WordNet (Fellbaum, 1998 and Princeton, 2010), a large lexical database of English, to extract the topological categorization. WordNet groups words into synsets, sets of cognitive synonyms, and models semantic relationships between them, one type of which is the *is-a* relation (for nouns), that forms a directed acyclic graph (DAG) with one root entity. For this reason, it is possible to traverse the hypernyms of the selected synsets to the root, and represent each word with the entity at a given depth in this path. If two words' least common subsumer (LCS) is of equal or higher depth (further from the root) than the chosen depth level, they will be represented by the same label. By varying the chosen distance from the root, we can tune the granularity of the categorization, e.g. the number of labels; we chose a depth of 2, resulting in 9 labels. See Figure 1 for an illustration.

4.4 Thesaurus-Affinity Based Grouping

We can extend this notion of similarity to calculate pairwise similarities between synsets using the Wu-Palmer Similarity measure (Wu and Palmer, 1994) and use the resulting similarity matrix as an input to a clustering algorithm. For consistency, we applied the affinity propagation algorithm (Section 4.2).

4.5 Out-of-Domain Clustering

The methods discussed so far build on the assumption that we have general information about target domains, that are otherwise not used to train the model. Namely, we build the categories either knowing the named entity labels occurring in such domains or create the categories such that it is likely that entity labels will have a sensible category to be grouped under.

This is a realistic scenario in cases where target domains are known in advance but no training data is available for them. However, to remove this assumption we assign unique entities for each domain (Appendix F.1) based on the occurrence count of the entity labels in different domains, except for the news domain that has only generic entity types.

For each domain (*evaluation domain*) we remove unique entity labels from the pool of all labels, cluster the remaining ones using the embedding space method described in Section 4.2, then predict which cluster the unique entity labels belong to. This results in one grouping per *evaluation domain* that will be applied only in the out-of-domain evaluation setting. (Section 5.2)

5 Experimental Setup

5.1 Model Implementation

To isolate the effect of our interventions, we base our implementation on the CrossRE project which follows the state-of-the art by Baldini Soares et al. (2019).

In the model, a sentence s containing an ordered pair of entities (e_1, e_2) is augmented by entity span markers e_1^{start} , e_1^{end} , e_2^{start} , and e_2^{end} . We enrich the delimiters (Zhong and Chen, 2021) with the cluster label of the entity. This is different from CrossRE, where the entity label itself is used. For example, following the replacements in the manual grouping magazine \rightarrow artifact and writer \rightarrow person:

<E1:artifact> The New Yorker
</E1:artifact> wrote that <E2:person>
Kubrick </E2:person> has taken [...]

As described in CrossRE, s is then fed into a pre-trained BERT encoder (Devlin et al., 2019) and the concatenated output representation of the entity start markers become the input to a single-layer feed-forward neural network with softmax activation function for the relation classification.

Since the models were trained on several domains at once, we implemented shuffling at training time to reduce bias and aid with generalization. The training parameters are specified in Appendix H.

5.2 Benchmark

We establish two baseline models that we propose to be used as benchmarks for further research. The first one combines cross-domain fine-tuning with *multi-domain evaluation*. In this case we train the model on the combined training dataset from all six domains and also test it on the combined test corpus. Therefore, it can be considered in-domain evaluation in a cross-domain setting. This setting models the few-shot learning scenario where all target domains are known in advance.

The second benchmark combines cross-domain training on five domains and *out-of-domain evaluation* on the sixth (test / target), done six times, resulting in one baseline for each domain. This mimics the use case where the model might be used on a previously unknown domain with unique entity types, or the target domain lacks training data.

5.3 Evaluation

As baseline models we use our benchmarks as described in Section 5.2 with a *no mapping* scenario, injecting the entity types themselves, not the clusters names. We compare these baselines to all combinations of mapping types and evaluation methods, except for the OOD Clustering, which is not relevant to the multi-domain evaluation setting.

When examining entity labels unique to a domain for OOD Clustering, we noted that some were present in the dataset of other domains. To realize the scenario for which we set up this mapping type (Section 4.5), we removed the sentences containing these entities from the corpus in this particular use case. For comparability, we preserved relevant training and test set sizes by randomly reallocating sentences from the dev set. This affected 0.44% of the total corpus, see Appendix F.3 for the detailed numbers.

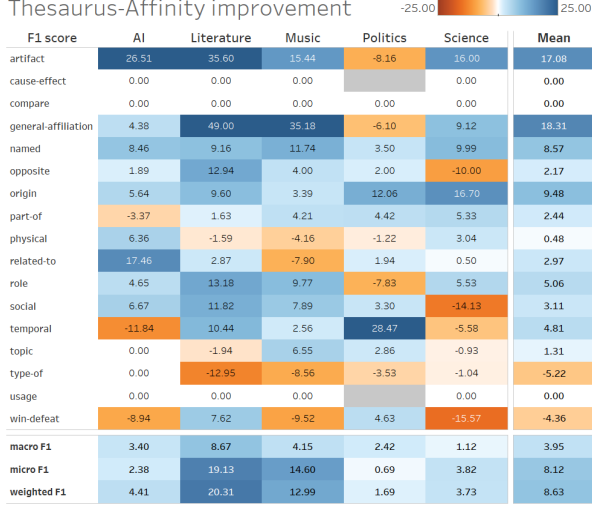


Figure 2: Difference in F1 scores for Thesaurus Affinity compared to baseline for each relation type and domain.

6 Results

6.1 OOD Evaluation Results

We chose to exclude the news domain from the OOD evaluation, due to the challenges described by Bassignana and Plank, 2022a, otherwise the highly imbalanced label distribution and many relation types having a support of zero would skew the results. Based on Harbecke et al., 2022 we used weighted F1 as our main metric, however the trends are equal for micro and macro F1. Detailed results can be found in Appendices I through J.

The best performing mapping methods are Manual and Thesaurus Affinity, outperforming the baseline by 7.6 and 8.6 percentage points, respectively. These methods perform significantly better than the baseline in every domain, as can be seen in Table 1, and all mapping methods perform better than the baseline overall. Figure 2 shows the F1 score difference for all relations in each domain for the Thesaurus Affinity compared to the baseline models in the OOD evaluation setting. Absolute results as heatmaps for all domains can be found in Appendix I.1, difference maps in Appendix I.2 and per domain performance charts in Appendix I.3.

From Figure 3 it can be seen that almost all mappings perform better than the baseline in each domain. The only exceptions are OOD Clustering in the AI domain as well as OOD Clustering & Embedding in politics. Specifics of the performance will be discussed further in section 7.1.

6.2 Multi-Domain Evaluation Results

The average weighted F1 scores for the Multi-Domain evaluation can be found in Table 1. It

is clear that the performance of the baseline is in the middle of the field, with the leaders Manual and Thesaurus Affinity not far ahead. A breakdown of F1 scores on relation label level can be found in Appendix J, and a summary table for both evaluation methods in Appendix K.

7 Discussion

7.1 OOD Evaluation

In general, all methods suffer from the label distribution not being uniform, as discussed by Bassignana and Plank, 2022a. This leads to relations with F1 scores of zero and makes it difficult to improve performance on these relation types. Future work might benefit from training with higher penalty on less frequent relation labels.

The **OOD Clustering** is the method with the least overlap of information between training and target domain, and thus also the hardest test case to beat. It is notable that it still performs significantly better overall than the baseline in most domains, thus it is the clearest indicator that entity type clustering generalizes better for out-of-domain use cases.

As expected, it performs worse than other grouping methods. This is partially due it being a more challenging subset of the **Embedding** clustering method, which itself performed only slightly better, leaving a significant performance gap compared to other mapping methods. Surprisingly, OOD Clustering performed better than Embedding in the science domain, which has the most unique entities, making it especially hard for the former. We leave detailed analysis of the connection between the number of unique entity types and the performance of OOD Clustering to future work. We also note affinity propagation’s tendency to group words by topic, contributing to its middling performance.

Manual mapping, Topological and Thesaurus Affinity methods are based on linguistics and how humans interpret language, which led to the best performance across domains. In **Topological**, WordNet’s hierarchical modeling of word senses is aligned with our aim to group entities based on semantic relations. However, commonality expressed through Wu-Palmer similarity appears to be an even better basis for clustering. Since **Thesaurus Affinity** uses the same clustering algorithm as Embedding, the performance difference shows that WordNet is able to capture semantic relations better than the word2vec model and the limiting

Eval. Type	Baseline	Manual	Embedding	OOD Clustering	Topological	Thesaurus Affinity
OOD eval.	43.54	51.15	47.94	46.22	51.03	52.16
Multi-Domain	61.86	62.19	61.15	-	61.38	61.98

Table 1: Weighted F1 scores for each mapping method and evaluation type.

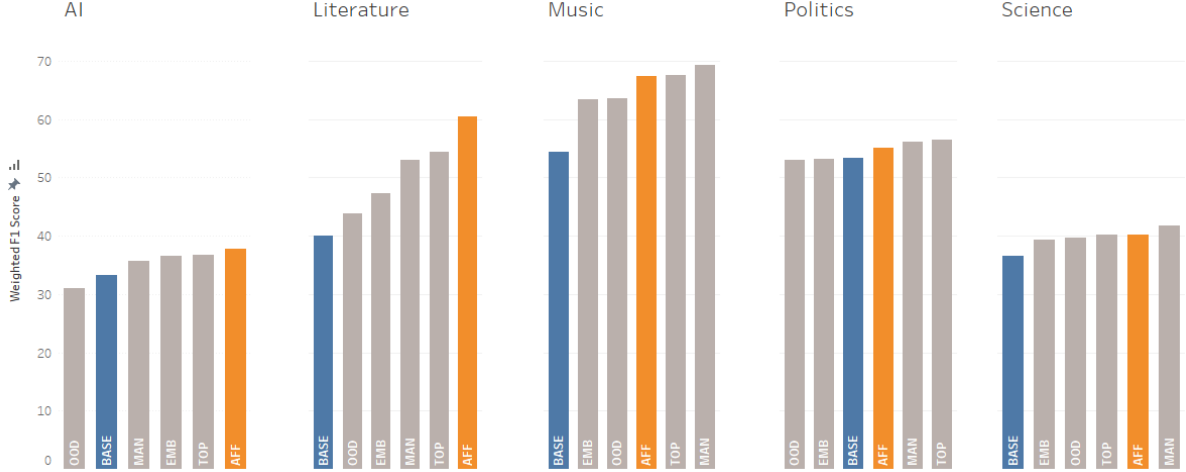


Figure 3: Weighted F1 scores for each mapping method over the five domains in OOD evaluation.

factor is not the algorithm.

In the **Manual** method we grouped entity types on the basis of our understanding of the labels and descriptions, leading to a performance nearly on par with the best result. We theorize however that this method is the most generalizable and with expert knowledge a mapping could be made which would suit numerous domains. Additionally, this does not suffer from the updatability issues of a thesaurus.

7.2 Multi-Domain Evaluation

As expected, all clustering methods perform better in this setting than their OOD counterparts. In Table 1 we can see that the performance is virtually the same across the board.

We find that the results, as displayed in Table 1, makes sense, as domain-specific entity types interact in a way that might not be caught by the less granular categories. Grouping entity types together in categories involves removing information and might therefore not capture the way entity types interact properly.

8 Conclusion

We researched how clustering named entity types in cross-domain relationship classification setups can increase performance and lead to more generalizable models.

We find that significant performance increases can be achieved by using entity mapping in such scenarios, where training data for the target domain is unavailable or the model is used for inference on a previously unseen domain. It is particularly relevant for use cases with unique entity types in the target domain.

On the other hand, it is not relevant in the case of target domain training data being available, as the granularity of the original entity types offer valuable domain-specific information.

We find that the combination of human centered word-understanding using WordNet and clustering (Thesaurus Affinity) was the best performing mapping method in the OOD evaluation setting, but that manually curated groups would create the most generalizable and sustainable mapping.

Limitations

The result are limited to the English language and the source data is mostly generated in the Global North by people with internet access.

The Manual grouping and simplifications are created by us, limiting us to our interpretation of the entities and descriptions.

Detailed analysis of metadata – syntax ambiguity and uncertainty – collected during the annotation process is outside of the scope of this project, leaving it for future work.

References

- Livio Baldini Soares, Nicholas FitzGerald, Jeffrey Ling, and Tom Kwiatkowski. 2019. [Matching the blanks: Distributional similarity for relation learning](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2895–2905, Florence, Italy. Association for Computational Linguistics.
- Elisa Bassignana and Barbara Plank. 2022a. [CrossRE: A cross-domain dataset for relation extraction](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3592–3604, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Elisa Bassignana and Barbara Plank. 2022b. [What do you mean by relation extraction? a survey on datasets and study on scientific relation classification](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*, pages 67–83, Dublin, Ireland. Association for Computational Linguistics.
- Emily M. Bender and Batya Friedman. 2018. [Data statements for natural language processing: Toward mitigating system bias and enabling better science](#). *Transactions of the Association for Computational Linguistics*, 6:587–604.
- 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.
- Christiane Fellbaum. 1998. [WordNet: An Electronic Lexical Database](#). Bradford Books.
- Brendan J. Frey and Delbert Dueck. 2007. [Clustering by passing messages between data points](#). *Science*, 315(5814):972–976.
- Tianyu Gao, Xu Han, Hao Zhu, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. 2019. [FewRel 2.0: Towards more challenging few-shot relation classification](#). 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 6250–6255, Hong Kong, China. Association for Computational Linguistics.
- Google. 2013. [word2vec - Google Code Archive](#). Last accessed 11 May 2023.
- David Harbecke, Yuxuan Chen, Leonhard Hennig, and Christoph Alt. 2022. [Why only micro-f1? class weighting of measures for relation classification](#). In *Proceedings of NLP Power! The First Workshop on Efficient Benchmarking in NLP*, pages 32–41, Dublin, Ireland. Association for Computational Linguistics.
- J.A. Lee and M. Verleysen. 2007. *Nonlinear Dimensionality Reduction*. Information Science and Statistics. Springer New York. Section 7.2.
- Zihan Liu, Yan Xu, Tiezheng Yu, Wenliang Dai, Ziwei Ji, Samuel Cahyawijaya, Andrea Madotto, and Pascale Fung. 2021. [Crossner: Evaluating cross-domain named entity recognition](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(15):13452–13460.
- Leland McInnes, John Healy, and James Melville. 2020. [UMAP: Uniform manifold approximation and projection for dimension reduction](#).
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. [Distributed representations of words and phrases and their compositionality](#).
- University Princeton. 2010. About wordnet. <https://wordnet.princeton.edu/>. Accessed: 2023-04-23.
- 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.
- Laurens van der Maaten and Geoffrey Hinton. 2008. [Visualizing data using t-SNE](#). *Journal of Machine Learning Research*, 9(86):2579–2605. Section 4.2.
- Zhibiao Wu and Martha Palmer. 1994. [Verb semantics and lexical selection](#). In *32nd Annual Meeting of the Association for Computational Linguistics*, pages 133–138, Las Cruces, New Mexico, USA. Association for Computational Linguistics.
- Zexuan Zhong and Danqi Chen. 2021. [A frustratingly easy approach for entity and relation extraction](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 50–61, Online. Association for Computational Linguistics.

Appendix

A Named Entity and Relation Types

A.1 Named Entity Labels

See Table 2 for the original labels as published in Liu et al. (2021) and labels added by Bassignana and Plank (2022a).

A.2 Semantic Relation Labels

As published in Bassignana and Plank (2022a).

Part-Of, Physical, Usage, Role, Social, General-Affiliation, Compare, Temporal, Artifact, Origin, Topic, Opposite,

Domain	CrossNER	CrossRE (added)
<i>News (Reuters)</i>	person, organization, location, miscellaneous	country
<i>Politics</i>	politician, person, organization, political party, event, election, country, location, miscellaneous	album, book, musical artist, musical instrument, song, university
<i>Natural Science</i>	scientist, person, university, organization, country, location, discipline, enzyme, protein, chemical compound, chemical element, event, astronomical object, academic journal, award, theory, miscellaneous	book
<i>Music</i>	music genre, song, band, album, musical artist, musical instrument, award, event, country, location, organization, person, miscellaneous	magazine
<i>Literature</i>	book, writer, award, poem, event, magazine, person, location, organization, country, miscellaneous	album, band, literary genre, music genre, programming language, scientist, song, university
<i>Artificial Intelligence</i>	field, task, product, algorithm, researcher, metrics, university, country, person, organization, location, miscellaneous	academic journal, conference, event, programming language

Table 2: Named Entity labels used in CrossNER (Liu et al., 2021) and added in CrossRE (Bassignana and Plank, 2022a) per domain.

Cause-Effect, Win-Defeat, Type-Of, Named, Related-To.

B Manual Groupings

B.1 Groups Obtained from Elisa Bassignana

person: person, researcher, musical artist, politician, scientist, writer

location: location, country

organisation: organisation, university, political party, band

event: event, election, conference

miscellaneous: miscellaneous, program language, metrics, product, field, poem, musical instrument, song, theory, academic journal, chemical element, protein, discipline, enzyme, chemical compound, astronomical object, award, music genre, album, book, magazine, task, algorithm, literary genre

B.2 Proposed Manual Grouping

person: person, researcher, musical artist, politician, scientist, writer

location: location, country

organisation: organisation, university, political party, band

event: event, election, conference, award

artifact: album, song, academic journal, poem, magazine, book

scientific: metrics, enzyme, protein, chemical compound, chemical element, astronomical object

concept: theory, music genre, field, discipline, algorithm, literary genre

brand: product, programming language

miscellaneous: miscellaneous, musical instrument, task

C Embedding Space Based Grouping

Category names were obtained automatically by finding the word closest to the mean of the entity labels in the category in the google-news-300 word2vec embedding space, excluding words containing punctuation and special characters. This categorization is an example, as the result depends on the random state.

essay: academic journal, book, magazine, poem, theory

songs: album, band, music genre, literary genre, musical instrument, song

proteins: algorithm, astronomical object, chemical compound, chemical element, enzyme, metrics, programming language, protein

convention: award, conference, election, event, location, organisation, person, political party

TELW: country, discipline, field, miscellaneous, product, task, university

journalist: musical artist, politician, researcher, scientist, writer

D Topological Grouping

Since one word can belong to several synsets, carrying different meanings, it is possible to disambiguate the meanings (e.g. java the programming language, not the island). We theorize that the effect of these substitutions is negligible because of the directed acyclic nature of the *is-a* relationships (Figure 1).

The category names are the synset identifier names on the selected level, with the number removed from the end (except if a name clash were to occur).

object: album, book, musical instrument, academic journal, location, magazine, product

psychological feature: algorithm, discipline, election, event, field, music genre, literary

communication: award, programming language, poem, song

group: band, conference, country, astronomical object, organisation, political party, university

matter: chemical compound, chemical element, enzyme

measure: metrics

assorted: miscellaneous

casual agent: musical artist, person, politician, researcher, scientist, writer

thing: protein

E Thesaurus-Affinity Based Grouping

Category names were obtained the same way as described in Appendix C. This categorization is an example, as the result depends on the random state.

journals: album, book, musical instrument, academic journal, magazine, product

pleiotropy: algorithm, discipline, election, event, field, music genre, literary genre, metrics, miscellaneous, task, theory

poems: award, programming language, poem, song

group: band, conference, country, astronomical object, organisation, political party, university

enzymes: chemical compound, chemical element, enzyme, protein

journalist: location, musical artist, person, politician, researcher, scientist, writer

F Out-of-Domain Clustering

F.1 Unique Entities

Artificial Intelligence: algorithm, conference, field, metrics, product, researcher, task, programming language

Literature: poem, writer, literary genre, book, magazine

Music: band, song, musical artist, music genre, album, musical instrument

Politics: political party, election, politician

Science: chemical compound, chemical element, astronomical object, discipline, enzyme, protein, theory, scientist, academic journal

F.2 OOD Clustering Example

An example for OOD clustering with Artificial Intelligence as evaluation domain. The predicted entities (unique to AI) are italicized. Category names were obtained the same way as described in Appendix C. This categorization is an example, as the result depends on the random state.

songs: album, band, music genre, literary genre, poem, song, musical instrument

group: award, political party, event, location, organisation, person, *conference, task*

novelist: book, election, academic journal, magazine, musical artist, politician, scientist, writer, *researcher*

proteins: chemical compound, chemical element, enzyme, astronomical object, miscellaneous, protein, *algorithm, programming language*

disciplines: country, discipline, theory, university, *field, metrics, product*

F.3 Sentences removed in OOD Clustering training

Domain	Removed	%
AI	3	0.34
Literature	10	1.09
Music	1	0.12
News	0	0
Politics	2	0.24
Science	7	0.82
Total	23	0.44

Table 3: Number and percentage of sentences removed in OOD Clustering training.

G Substitutions for clustering

Named entity label substitutions. Substitutions specific to Embedding Space Based Grouping and Out-of-Domain Clustering are marked with word2vec, while those specific to Topological Grouping and Thesaurus-Affinity Based Grouping with WordNet.

Original	Substitute
musical artist	musician
organisation	organization
political party	coalition (word2vec), party (WordNet)
academic journal	journal
chemical compound	chemical
chemical element	chemical
astronomical object	galaxy
music genre	genre
literary genre	genre
programming language	javascript (word2vec), java (WordNet)
musical instrument	violin (word2vec), instrument (WordNet)
misc	miscellaneous

Table 4: Named entity label substitutions.

H Reproducibility

The hyperparameters are reported in Table 5. They are used across all mapping types & the baseline, and across both evaluation types. The early stopping parameter stops the training if there are no improvements on the development data for three consecutive epochs. The training were divided between two machines with the following specifications:

- Nvidia® GeForce™ RTX 3090 Ti 24 GB GPU; Intel® Core i9-9900K CPU
- Nvidia® GeForce™ GTX 1660 Super 6 GB GPU; Intel® Core i5-11400F CPU

Parameter	Value
General	
Random Seeds	4012, 5096, 8857, 8878, 9908
Clustering	
PCA n_components	35 ⁵
UMAP n_components	3
distance metric	cosine
UMAP n_neighbor	4
UMAP min_dist	0.3
Affinity Prop. damping factor	0.5
Model	
Encode	bert-base-cased
Classifier	1-layer FFNN
Loss	Cross Entropy
Optimizer	Adam
Epochs	50
Batch Size	32
Learning Rate	$2e^{-5}$
Early Stop	True
Shuffle Data	True

Table 5: Hyperparameter settings.

I OOD Evaluation Result Figures

I.1 Heatmap over relation and domain

Gray colored fields indicate a case where the support is zero for a relation in the test set.

⁵When the number of entity types in the OOD Clustering case is below 35, then that number is used.

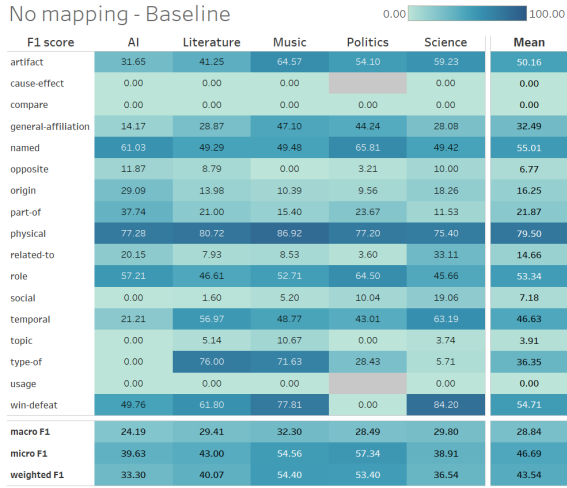


Figure 4: Baseline Heatmap - Average F1 score

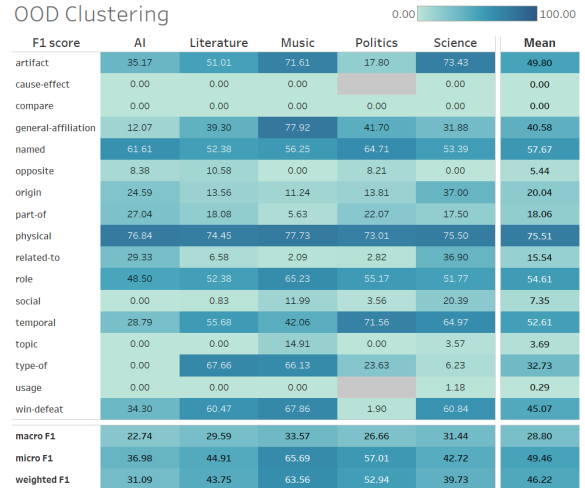


Figure 7: OOD Clustering Heatmap - Average F1 score

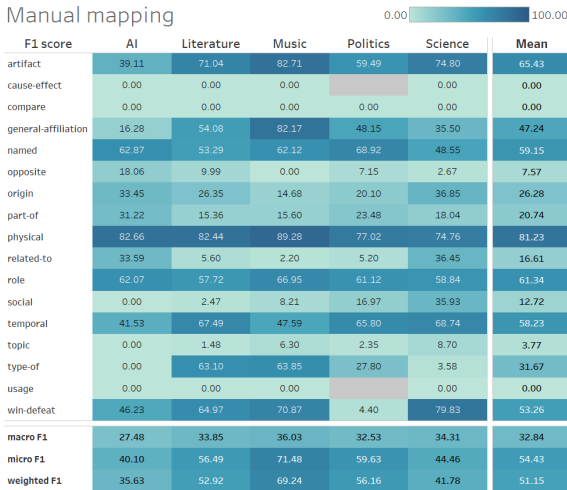


Figure 5: Manual Heatmap - Average F1 score

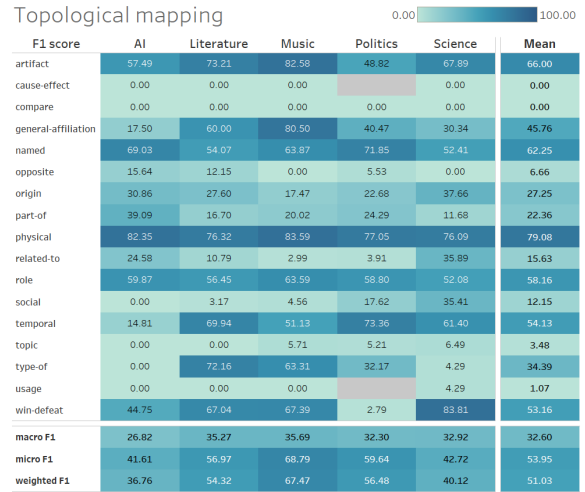


Figure 8: Topological Heatmap - Average F1 score

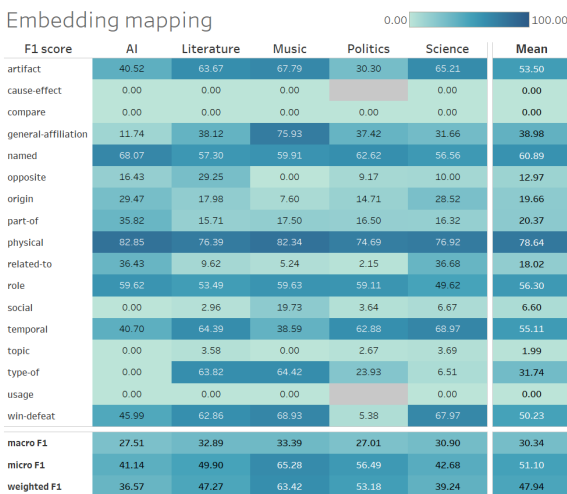


Figure 6: Embedding Heatmap - Average F1 score

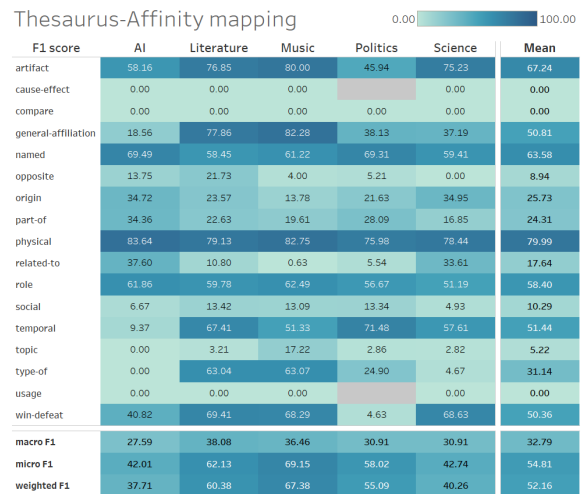


Figure 9: Thesaurus Affinity Heatmap - Average F1 score

I.2 Heatmap over relation and domain - difference from baseline

Gray colored fields indicate a case where the support is zero for a relation in the test set.

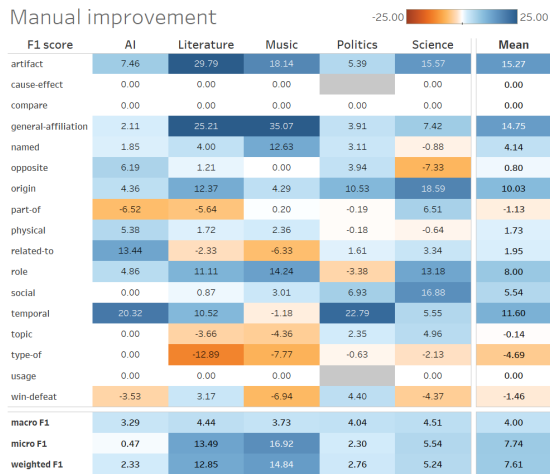


Figure 10: Manual Heatmap - difference from baseline

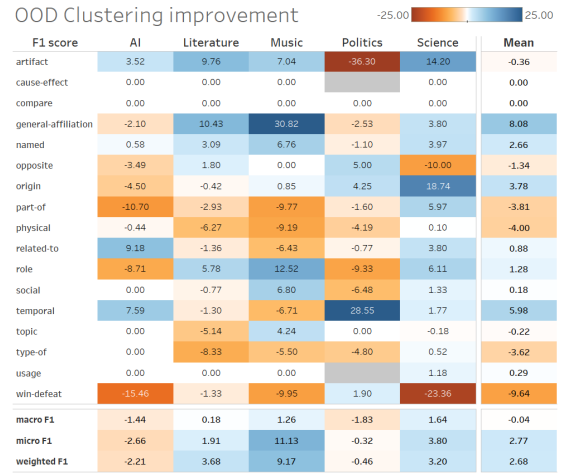


Figure 12: OOD Clustering Heatmap - difference from baseline

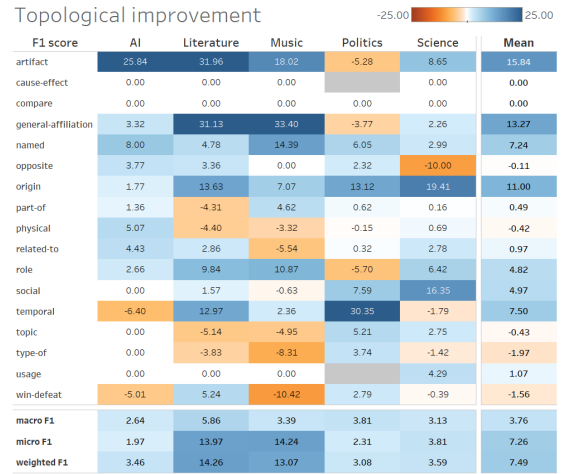


Figure 13: Topological Heatmap - difference from baseline

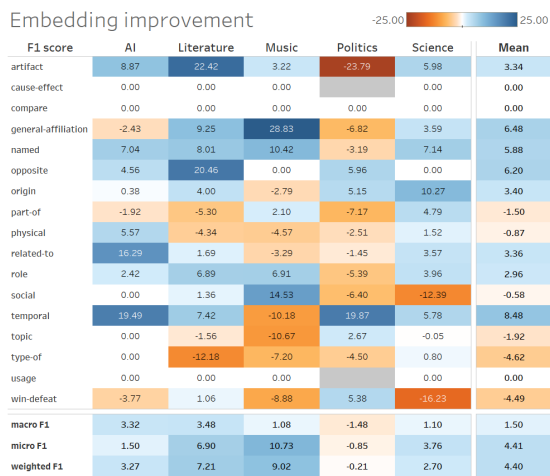


Figure 11: Embedding Heatmap - difference from baseline

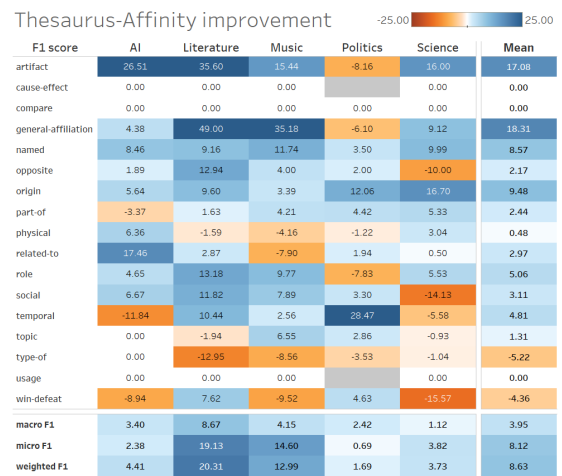


Figure 14: Thesaurus Affinity Heatmap - difference from baseline

I.3 Micro & Macro F1 Results Across Domains

See Figure 15 for macro and Figure 16 for micro F1 scores for each mapping method over the five domains in OOD evaluation. The Baseline (no mapping) is colored blue, and the overall best model, Thesaurus Affinity, is colored orange.

J Multi-Domain Result Figures

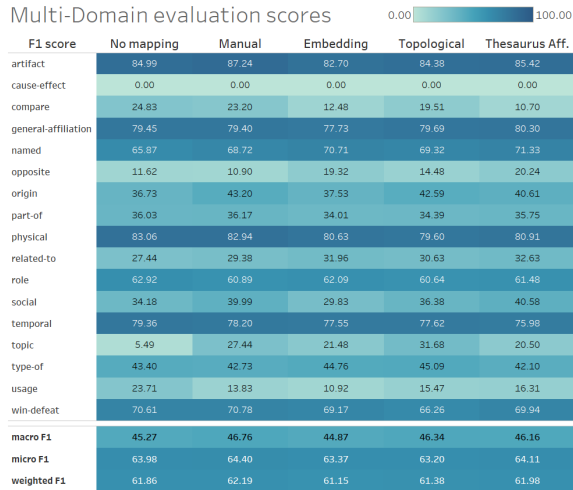


Figure 17: Heatmap over mapping-type and relations - F1 score

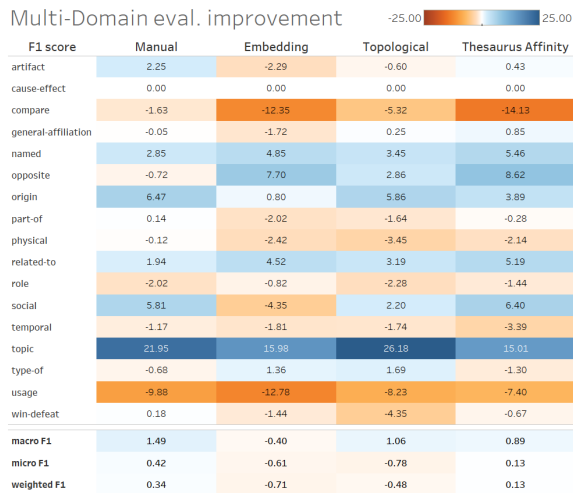


Figure 18: Heatmap over mapping-type and relations - F1 score gain over Baseline

K Results Summary

A summary of micro, macro, and weighted F1 scores for each mapping method and evaluation type can be found in Table 6.

L Group Contributions

The overall workload was equally distributed across group members. Throughout the project we have worked both in conjunction and individually, thus all parts of the project has been done and/or reviewed by all group members.

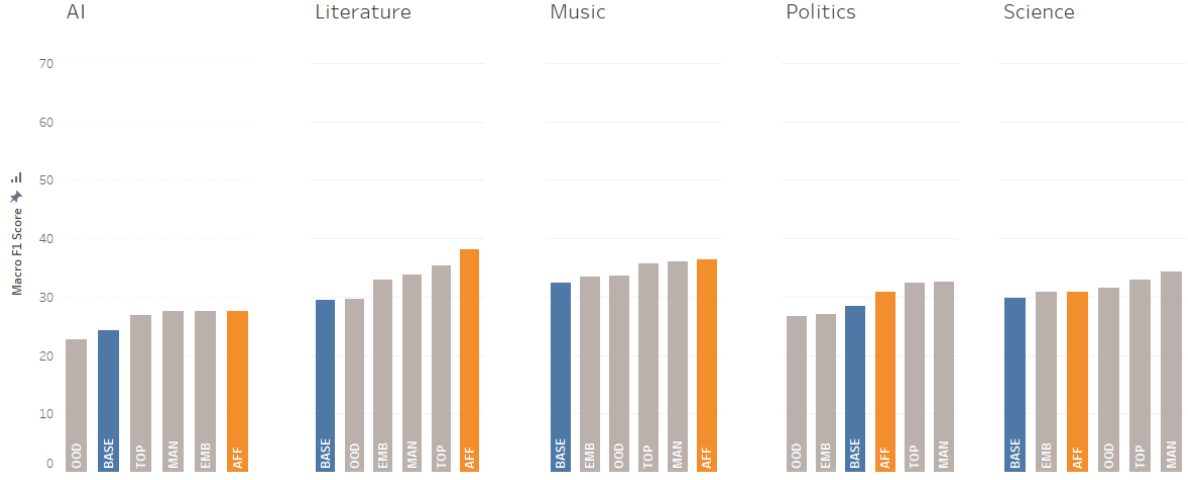


Figure 15: **Macro F1** scores for each mapping method over the five domains in OOD evaluation.

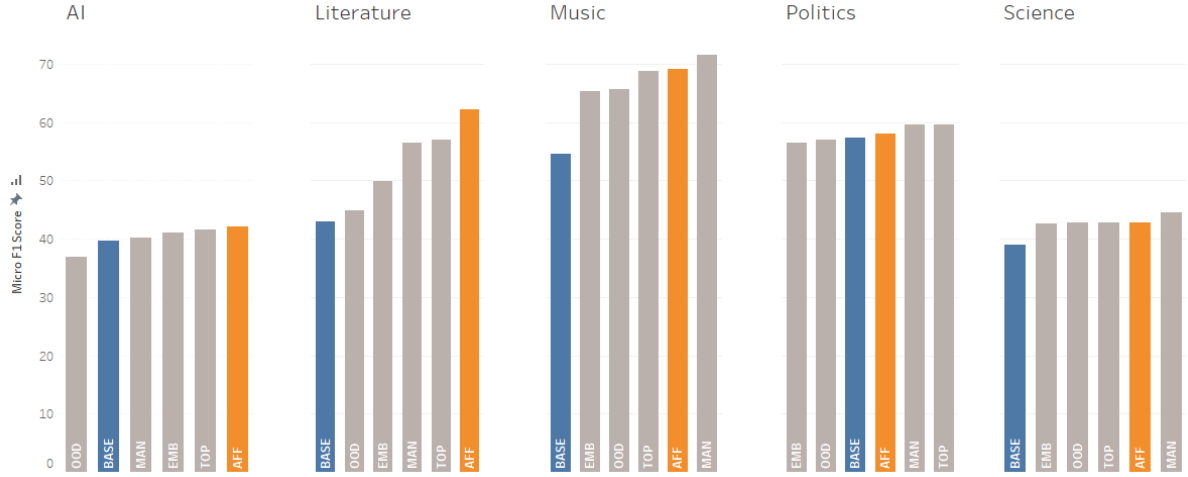


Figure 16: **Micro F1** scores for each mapping method over the five domains in OOD evaluation.

Eval. Type	Metric	Baseline	Manual	Embedding	OOD Clustering	Topological	Thesaurus Affinity
OOD	Weighted F1	43.54	51.15	47.94	46.22	51.03	52.16
	Macro F1	28.84	32.84	30.34	28.80	32.60	32.79
	Micro F1	46.69	54.43	51.10	49.46	53.95	54.81
Multi-Domain	Weighted F1	61.86	62.19	61.15	-	61.38	61.98
	Macro F1	45.27	46.76	44.87	-	46.34	46.16
	Micro F1	63.98	64.40	63.37	-	63.20	64.11

Table 6: F1 scores for each mapping method and evaluation type.