

Injecting Explicit Cross-lingual Embeddings into Pre-trained Multilingual Models for Code-Switching Detection

Anonymous ACL submission

Abstract

Code-switching has become the modus operandi of internet communication in many communities such as South Africans, who are domestically multilingual. This phenomenon has made processing textual data increasingly complex due to non-standard ways of writing, spontaneous word replacements, and other challenges. Pre-trained multilingual models have shown elevated text processing capabilities in various similar downstream tasks such as language identification, dialect detection, and language family discrimination. In this study, we extensively investigate the use of pre-trained multilingual models - AfroXLMR, and Serengeti for code-switching detection on five South African languages: Sesotho, Setswana, IsiZulu, IsiXhosa, and English, with English used interchangeably with the other four languages, including various transfer learning settings. Additionally, we explore the modelling of known switching pairs within a dataset through explicit cross-lingual embeddings extracted using projection models: VecMap, Muse, and Canonical Correlation Analyses (CCA). The resulting cross-lingual embeddings are used to replace the embedding layer of a pre-trained multilingual model without additional training. Concretely, our results show that performance gains can be realized by closing the representational gap between the languages of the code-switched dataset with known codes, using cross-lingual representations. Moreover, expanding code-switched datasets with datasets of closely related languages improves code-switching classification, especially in cases with minimal training examples.

1 Introduction

Code-switching, a term used interchangeably with code-mixing in Natural Language Processing (NLP), refers to a linguistic phenomenon, where a single utterance or text is made up of multilingual

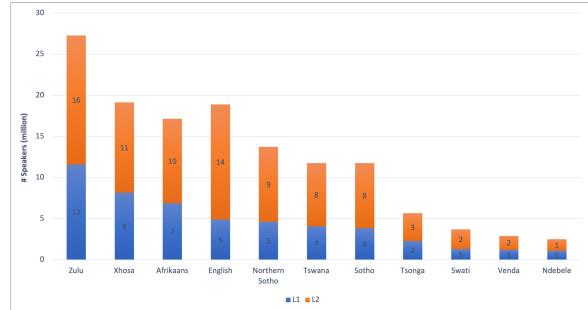


Figure 1: Number of First Additional (L1), and Second additional Speakers (L2) in a South African Household.

tokens (referred to as codes), arranged meaningfully to the receiving audience, prevalent in multilingual communities such as the Internet (Jose et al., 2020). In a typical South African household setting, the majority of families are multilingual as seen in Figure 1¹. However, language technologies developed for these South African low-resourced languages are monolingual in nature. This myopic view limits the ability of these users to fully express themselves within these technologies, thus, motivating a need for modelling code-switching in speech or text understanding technologies.

Code-switching is not a new challenge. In 2013, Modipa et al. (2013) investigated the implication of code-switching for developing automated speech recognition systems using Sepedi and English code-switched datasets sourced from radio broadcasts. Vyas et al. (2014), investigated code-switching detection formalized as language detection together with Part-of-speech (POS) tagging. While Singh et al. (2018) explored the creation of POS datasets from code-switched tweets that closely resemble real-world scenarios. On the other hand, Aguilar et al. (2019) proposed a shared task for modelling Named Entity Recognition (NER) on code-switched datasets. Closer to our methodol-

¹https://en.wikipedia.org/wiki/Languages_of_South_Africa

ogy, Hussain and Arshad (2021) proposed a hybrid model composed of an attention mechanism and recurrent neural network for code-switching detection.

However, none of the aforementioned works attempted a natural standpoint that brings the involved languages in the code-switched dataset closer together through cross-lingual representation learning. Cross-lingual representation learning emerges from a broader field of cross-lingual models, where the idea is to fuse two or more syntactic and semantically sound monolingual embeddings in order to supplement the linguistic shortcomings of both embedding spaces for improving downstream task performance (Mikolov et al., 2013b; Faruqui and Dyer, 2014; Artetxe et al., 2018). We hypothesize that given known languages $L_g = \{l_1, l_2, \dots, l_n\}$ in the labelled code-switched dataset C_{cs} , developing a language model M_L through cross-lingual representation learning of languages in L_g can improve on the task of code-switching classification (and possibly generalize to other extensions of code-switching tasks such as code-switched POS, code-switched NER, and more complex tasks like code-switched question answering) following the linguistic equivalence constraint of code-switching. Using cross-lingual representation learning to bring the languages in the code-switched datasets closer together, may establish a representational space that closely matches code-switching, thus, improving training, contrary to treating the codes as separate entities during training. Subsequently, we then want to answer the question - what effect do explicit cross-lingual representations, generated from the involved language, have on the performance of code-switching classification? To answer this, we conduct a systematic analysis and evaluation of supervised models using both monolingual and cross-lingual embeddings on the task of word-level code-switching classification. The field of cross-lingual models has witnessed great success over the years in NLP and has evolved to advanced processing capabilities using large pre-trained multilingual models (Devlin et al., 2018; Conneau et al., 2019). On this, our methodology investigates the downstream performance of code-switching classification prior to the injection of cross-lingual embeddings into the embedding layer of transformer architectures AfroXLM-r, and Serengeti as our training and classification models.

In summary, our major contributions are suc-

cinctly organized as follows:

- We investigate the use of cross-lingual embeddings in the context of code-switching classification for four South African low-resource languages, namely, Setswana, Sesotho, isiZulu, and IsiXhosa, and one high-resourced language - English.
- We showcase the experimental results that highlight an intuitive relationship between code-switching and cross-lingual representations.
- We highlight characteristics that are deterministic of improved performance for code-switching classification using cross-lingual and monolingual embeddings.

2 Related Work

Recent studies show the need to accelerate code-switching inclusion in language technology developments and the limitations incurred if otherwise (Sitaram et al., 2019). In line with this shared responsibility, various works addressing different challenges across different downstream tasks have been proposed for code-switching. For code-switching detection, an early work by (Chittaranjan et al., 2014) proposed a Conditional Random Field (CRF) model based on their success in sequence labelling and a subword segmental approach that segments lexical units according to their morphological grounding. Xia (2016) explores subword vectors from (Bojanowski et al., 2017), prefix and suffix extraction, and linear-chain CRF for code-switching detection and reports (at the time) promising results (83.0 %, 94.9 % F1 score, and accuracy respectively). While Patro et al. (2017), remedies incorrect annotations of borrowed words in the code-switched text by proposing a set of likeliness metrics that utilizes language usage patterns on Twitter.

Code-switching classification as a stand-alone downstream task has use-case limitations, and as such various works have coupled it with other downstream tasks. Solorio and Liu (2008) exploits existing probabilistic tree-based taggers to generate training features of the switched languages for POS tagging in the context of code-switching. They use heuristic-driven methods to combine these information-rich features to assist machine-learning-based prediction. Barman et al. (2016)

uses a joint Factorial CRF model to process complex trilingual code-mixed data for code-switching detection and POS tagging simultaneously. In contrast, Van der Goot and Çetinoğlu (2020) proposes code-switching-tailored lexical normalizers to improve the impact of non-canonical data points on POS tagging.

For code-switching in the context of NER, Attia et al. (2018) takes an architecture-driven approach and proposes a hybrid model consisting of enriched pre-trained embeddings, a Bidirectional Long and Short Term Memory (Bi-LSTM) model to the left and right context over the continuous representations, a convolution layer to model spatial dependencies, and finally a CRF for sequence prediction. Closer to our work, Winata et al. (2018) concatenates English and Spanish monolingual pre-trained embeddings together with character-level representation to address the out-of-vocabulary (OOV) issue. In contrast to our work, they did not explore cross-lingual representations. In the advent of the attention era, Winata et al. (2019) addresses both NER and code-switching detection using multilingual meta-embeddings extracted using a fully connected neural network and an attention mechanism layer.

In the aforementioned works, South African low-resourced languages lag behind across multiple tasks and techniques for processing code-switching text. In an attempt to remedy this limitation as well as ignite research interest, we propose a systematic analysis of using both monolingual and cross-lingual vector representations for code-switching detection together with injecting static cross-lingual embeddings into the embeddings matrix of the transformer architectures.

3 Methodology

3.1 Corpora

This study explores four low-resourced South African languages (LRSAL), namely, Sesotho, Setswana, isiZulu, and IsiXhosa, and one high-resourced language - English with codes: sot, tsn, zul, xho, and eng respectively. The sources for our monolingual corpora are: Flores NLLB Team (2022), WMT Costa-jussà et al. (2022), MC4 Raffel et al. (2019), NCHLT (Eiselen and Puttkammer, 2014), and African Crawl Dataset Végi et al. (2022). The statistics of our monolingual data are outlined in Table 1

Lang.	Sentences	Unique Voc	After LID	Unique Voc
tsn	1.1M	388K	462K	118K
sot	2.6M	1.5M	750K	453K
xho	2.7M	2.8M	1.2M	1.2M
zul	7.6M	9.2M	1.5M	1.5M

Table 1: Monolingual dataset description.

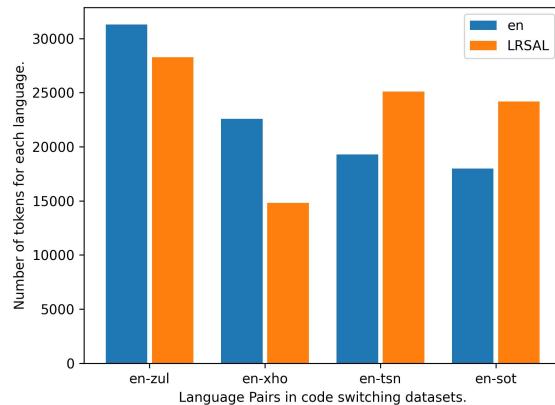


Figure 2: Code distributions in code-switched datasets.

3.2 Bilingual Lexicons

We used bilingual lexicons to generate the cross-lingual representation from the following sources: CPUT², Open Education Resource Term Bank (OERTB) (of Pretoria), and our manually collected lexicons from government public school repositories³. The collection resulted in 8742 en-tsn, 8763 en-sot, 11117 en-xho, and 17406 en-zul bilingual lexicon pairs.

3.3 Code-Switched dataset

We used the labeled code-switching dataset sourced from South African soap operas (Niesler et al., 2018). This dataset covers a continuum of spontaneous code-switching types such as intersentential, intrasentential, insertion, and etc. The languages covered in this dataset are as follows: English-isiZulu, English-isixhosa, English-Sesotho, and English-Setswana. Figure 2 shows the switching distribution of each pair in the datasets.

3.4 Pre-processing

Manual inspection of individual sources of monolingual corpora indicated that sources such as WMT, MC4, etc. contain foreign sentences and warranted further cleaning. For this, we used a publicly available language identification tool -

²<https://mlg.cput.ac.za/>

³Anonymous

242 GlotLID (Kargaran et al., 2023), which has shown
243 to have superior language identification perfor-
244 mance in Sindane and Marivate (2024). Since
245 the monolingual corpus is arranged by language,
246 we used GlotLID to confirm if the sentence s_i of
247 a known corpus C_{lang} belongs to language $lang$.
248 If GlotLID identifies s_i as not of $lang$ of corpus
249 C_{lang} it is discarded. Table 1, shows the remaining
250 number of lines before and after GlotLID. Further
251 pre-processing included the removal of URLs,
252 numbers, punctuations, and then lowercasing all words.

253 3.5 Monolingual Embeddings

254 Monolingual embeddings are continuous vector
255 representations of words (Mikolov et al., 2013a).
256 Various studies for generating these embeddings
257 have been proposed, with the recent FastText tech-
258 nique (Bojanowski et al., 2017) showing improve-
259 ment over previous methods such as Word2Vec
260 ⁴, and GloVe (Pennington et al., 2014) for low-
261 resourced settings. As such we generated mono-
262 lingual embeddings using FastText with dimension
263 $d = \{50, 100, 150, 200\}$. For English, we used
264 the available largely pre-trained embeddings from
265 GloVe, due to having the desired 3 dimensions d
266 $= \{50, 100, 200\}$. To evaluate the intrinsic quality
267 of all our embeddings we used the SimLex-999
268 dataset released by (Makgatho et al., 2021) con-
269 sisting of Setswana and Sepedi paired words to-
270 gether with their English translations. From this,
271 we derived that monolingual embeddings for low-
272 resourced languages require more data to train as
273 reflected by their inability to capture similar words
274 with similar representations as compared to En-
275 glish representations (Appendix B.1). However,
276 projecting the inefficiently learned representations
277 into a shared space with an effectively learned En-
278 glish representation space using the 3 techniques
279 (Canonical Correlation Analyses (CCA) (Faruqui
280 and Dyer, 2014), VecMap (Artetxe et al., 2018),
281 and MUSE (Lample et al., 2017)) shows improved
282 word similarity measures between intra-similar
283 (within a language) and inter-similar (between lan-
284 guages) (Appendix B.2). We further calculated
285 the Spearman’s correlation of the datasets and this
286 shows that cross-lingual embeddings outperform
287 monolingual embeddings. An extension of this
288 analysis is available in Appendix B.3, showing
289 plots for all embedding dimensions and the 3 pro-
290 jection types to evaluate the intrinsic quality of the

embeddings.

Notably, lower embedding dimensions of unre-
291 lated languages (xho and zul) show higher cosine
293 scores for paired word similarity evaluation.
294

295 3.6 Cross-lingual Embeddings

296 Cross-lingual embeddings are shared vector repre-
297 sentations generated using projection techniques
298 that aim to join monolingual embeddings of two
299 or more languages together with the objective of
300 transferring common desirable linguistic proper-
301 ties (Mikolov et al., 2013b). Projection techniques
302 develop a mathematical model driven by available
303 supervision resources such as bilingual lexicons,
304 parallel sentences, objective functions, etc. that
305 aims to learn how to translate source embeddings
306 (typically of high-resourced language) to target em-
307 beddings (of low-resourced) by transforming (shift,
308 distort, etc.) the source embeddings into a shared
309 vector space. In this study, we compared three
310 pioneering projection techniques – Canonical Cor-
311 relation Analyses (CCA), VecMap, and MUSE to
312 generate cross-lingual shared representations.

313 3.7 Pre-trained Models

314 Pre-trained multilingual models such as
315 mBERT (Devlin et al., 2018), RemBERT,
316 XLM-r (Conneau et al., 2019), and their Afro-
317 centric counterparts: Afri-BERTa (Ogueji et al.,
318 2021), Afro-XLM-r (Alabi et al., 2022) have
319 shown astonishing results for many downstream
320 tasks such as NER, POS Tagging, Machine
321 Translation, etc., with Afro-centric methods
322 having a slight performance edge over massively
323 pre-trained multilingual models with minimal to
324 no exposure to African languages. As such, we
325 will only concentrate on the Afro-centric model
326 Afro-XLM-r (base and large), and the recent
327 model Serengeti (Adebara et al., 2023) due to their
328 performance gains on downstream tasks.

329 3.8 Cross-lingual Injected Pre-trained Models

330 Suppose we know the languages $L_g =$
331 $\{l_1, l_2, \dots, l_n\}$ of the code-switched corpus
332 C_{cs} . We know that the interchangeable use of
333 linguistic tokens t_i, \dots, t_k from $\{l_1, \dots, l_n\}$ is not
334 random but rather orchestrated meaningfully to
335 bring the targeted languages coherently together.
336 Therefore, we hypothesize that creating shared
337 cross-lingual representations of the known
338 target languages that meaningfully tie the target
339 languages semantically together may improve

⁴<https://code.google.com/p/word2vec/>

code-switching processing. Concretely, this study’s code-switched datasets are a switching condition between *eng* and the four languages *sot*, *tsn*, *xho*, and *zul*. Hence, for known pairs (e.g *eng – sot*), we could create shared representations between *eng*, and *sot* into a meaningful shared vector space, such that translation pairs are semantically connected (i.e translation pairs between the two languages having similar representations), which could reflect in better processing of the target code-switching dataset. To evaluate this theory, we devised a technique that explores the use of cross-lingual embeddings into transformer architectures by replacing the embeddings layer with the new shared representation of the known pairs. Monolingual embeddings (to model an instance where the languages are not brought closer), and cross-lingual embeddings (to model a case where languages are brought closer together) will be used to replace the embedding layer of the transformer architectures. Experiments of these two setups will be compared and contrasted.

3.8.1 Dimensionality Reduction and Expansion

Our cross-lingual embeddings dimensions of $d = \{50, 100, 150, 200\}$ are significantly lower than the embeddings of AfroXLMr, and Serengeti of 700+. Therefore, we are tasked with either reducing the embeddings of the transformer when combined with cross-lingual embeddings or expanding cross-lingual embeddings to match existing transformer embeddings. In this case, we experimented with randomly initialized paddings with normal distribution to expand the cross-lingual embeddings to match transformer dimensions. For dimensionality reduction, Principal Component Analyses (PCA), Uniform Manifold Approximation and Projection, or other techniques can be adopted to reduce the existing transformer embeddings to match the static embeddings shape. Due to limited space, we leave this for future work.

3.9 Experimental Design

3.9.1 Corpora sizes

Following the pre-processing step, we extracted 80% of the most frequent words to generate our word embeddings.

3.9.2 Code-Switched dataset

We explored various setups to divide the dataset into train, development, and test sets. The con-

ventional setup included having a uniform train set across all languages of 4242 sentences and a test size of 1000 sentences, with only varying development sizes, since the size of the datasets for each language was not equal. The second setup included combining the dataset of the previous step to explore transfer learning. We combined code-switched datasets containing closely related languages (*xho* and *zul*) – two Nguni languages, and (*sot* and *tsn*) – two Sotho-Tswana languages. Finally, for each group, we added, a new language coming from a different language family, and then our last setup combined all languages.

3.9.3 Monolingual Embeddings generation

Our monolingual embeddings were trained for 50 epochs, with mostly default set-ups for FastText except for min character and maximum character considerations of 1, and 5 respectively. All words were reduced to lowercase.

3.9.4 Cross-lingual Embeddings generation

All projection techniques CCA, VecMap, and MUSE use a supervised projection setup with all available lexicons in this study. We used the default hyper-parameter set-up recommended by each technique’s proposed paper since no available resources exist to evaluate the intrinsic quality of cross-lingual representations.

3.10 Pre-trained Models

Each pre-trained model was trained for 20 epochs, used a batch size of 16, a maximum sequence cut-off of 200, a learning rate of $5 \exp^{-5}$ following Adelani et al. (2021), and Dione et al. (2023). However, Afro-XLMr-large did not perform well for the setup, and its hyperparameters were changed to a learning rate of $2 \exp^{-5}$, a batch size of 32, and was trained for 10 epochs.

4 Results

This section discusses the results of this study’s experimental findings. We report F1 score instead of accuracy since there is clear class imbalance between the codes within our datasets (Figure 2).

4.1 Baselines

Table 2, reports the F1 score performance of our baselines models: Afro-XLM-r base (b), Afro-XLM-r large (l), and Serengeti. Our results show that the models perform on par for word-based code-switching detection. Notably, increasing the

dataset size by combining datasets (e.g combining engxho with engzul to make engxhozul) shows to have a contradicting impact on code-switching detection. Firstly, combining datasets shows a drop in performance compared to monolingual training. Secondly, training a three pair of either 2 Nguni and one Sotho-Tswana or vice versa shows a 1-point drop when training with two Nguni (*xho*, *zul*), and one Sotho-Tswana (*tsn*), compared to training with only the two Nguni languages. We hypothesize that the addition of a language from a different language family creates an interference in the internal representation largely skewed to Nguni, which negatively impacts the model’s performance. Conversely, the Sotho-Tswana and single Nguni trio illustrate a performance gain in this setup. Since this pattern occurs on all models, it may imply that the internal structures of Sotho-Tswana may be robust and less susceptible to interferences compared to Nguni language learning. However, we leave this investigation for future work and continue investigating transfer learning.

4.2 Baselines Transfer learning

We investigate three modes of transfer learning in this section. The first setup investigates transferring models trained with multiple combinations above onto original (non-combined) datasets. The second transfer setup investigates language family grouping, where languages belonging to the same family (e.g *xho* and *zul*) are grouped into one label *nguni*, and *sttn* for Sotho-Tswana (results reported Table 3). Surprisingly, Afro-XLMR-b performs better than Serengeti on average for family-based-grouping code-switching detection. Additionally, all 3 models show higher performance for *nguni* combinations over Sotho-Tswana combination datasets. In the last setup, however, we changed the datasets to only have two labels (*eng*, and *swtc*), by replacing any other label that is not *eng* to *swtc*. This is done, to investigate if transfer could be easily modeled if the label set is reduced in a multi-code switching dataset. The results for the aforementioned last experiments are reported in Table 4. From the onset, the results show performance gains on language combinations compared to baseline results in Table 2. We hypothesize that this happens because the models may find it easy to create two representational clusters of the code-switching as opposed to creating multiple for each language. Additionally, closely related languages may not

cause incorrect predictions as in the conventional setup, as these are viewed as the same class *swtc*. Which is also supported by higher scores for original dataset scores. Regardless, this approach may be a better alternative to code-switching detection as many African languages are not effectively supported at the onset (i.e. pre-training), where, instead we model code-switching as certainty (*eng*) and uncertainty (not *eng*) binary classification with the assumption that English will easily be detected with high confidence. With this setup, smaller models such as Naive Bayes, SVM, etc, can be used for language identification of the *swtc* tag for finer-grained detail extraction. However, this claim requires additional empirical evidence and we leave it for future works. Finally, on average, Serengeti performs better compared to the two variations of Afro-XLM-r for this transfer category.

4.3 Cross-lingual Injected Pre-trained Models

Table 5, shows the experimental results of this study using various cross-lingual embeddings settings. Due to limited space, we only reported the results of Afro-XLM-r-base. We presented the remaining results in Appendix C with accompanying discussions and the Serengeti model results. From these results, CCA and VecMap embeddings show better performance for lower dimensions 50 and 100 over Muse embeddings, while Muse surpasses these two techniques on the highest dimension 200 on the original datasets *xho*, *zul*, *sot*, and *tsn*. This behavior is not clear as to why it occurs as we expected Muse embeddings to perform better than the other two techniques as shown by its representation quality illustrated in the intrinsic evaluation (Appendix B.2, Appendix B.3). This could mean that intrinsic performance is not correlated to extrinsic performance for this task. On the combination datasets, only Muse embeddings show consistent performance while the performance of CCA and VecMap fluctuates depending on the data combinations. This could be due to that, while the English to low-resourced language (*en* – LRL_i) pair remains semantically connected through the projection, the inclusion of a foreign pair *en* – LRL_k in the same space introduces the same issue (i.e the disconnect) we are trying to solve in the embedding space, for i , and k being low-resource languages. For example, combining en-*tsn*, with en-*xho* datasets, as these were not explicitly connected through projection techniques.

Baselines	Code-Switching dataset with base (eng)										
	xho	zul	sot	tsn	xhozul	sottsn	xhozulsot	xhozultsn	sottsnxho	sottsnzul	xhozulsottsn
Afro-xlmr-b	98	97	91	93	92	83	92	92	87	87	87
Afro-xlmr-l	97	97	90	90	91	82	91	92	85	86	86
Serengeti	97	97	91	92	92	83	91	92	87	87	87

Table 2: Reports the baseline model’s F1-score for code-switching detection for each dataset averaged over 5 runs.

Transfer	Code-Switching dataset with base (eng) and labels ngun and sttsn						
	xhozul	sottsn	xhozulsot	xhozultsn	sottsnxho	sottsnzul	xhozulsottsn
Afro-xlmr-b	98	93	95	96	94	94	95
Afro-xlmr-l	97	92	95	96	94	94	95
Serengeti	97	92	95	96	94	93	94

Table 3: Reports model’s F1-score for code-switching detection for each dataset averaged over 5 runs of grouped label datasets. The *xho*, *zul* are grouped into Nguni, and *sot tsn* grouped into Sotho-Tswana (*sttsn*).

To understand the impact of these injected cross-lingual embeddings we trained the best-performing model - Afro-XLMr-base, with monolingual embeddings from which these joined representations were formed. Table 16, shows these results for non-projected embeddings of *xho*, *zul*, *sot*, *tsn* created using FastText. From this, we can see that monolingual performance is significantly lower for all of the various evaluation metrics Accuracy, Precision, Recall, and F1-score. This in a way justifies our hypothesis that, perhaps not treating the codes as independent spaces but rather semantically joining their individual spaces at pre-training could improve code-switching processing as these languages are brought to the same shared space.

5 Analyses and Discussion

We also aimed to investigate the extent to which shared representations support text processing, especially for unknown, or out-of-vocabulary words for improved performance. That is, in cases where vocabulary is missing, we want to analyze if other related words are used for improving performance and if this is traced to shared representations. We aimed to achieve this through embedding attention score analyses. This is done by plotting embedding attention scores of predicted sentences to see where priority is placed for certain words.

Sentence attention score analyses consider relationships of words within a sentence, we also wanted to consider alternative words within the entire embedding matrix that may be important in processing the final output of the layer. We observed that high attention was also given to words outside the main sentence vocabulary. This could mean that, related words, were identified (possibly made possible by shared representations), thus, improv-

ing the processing of the target text. Monolingual plots do not show this behavior of exterior attention. Which implies that, indeed deeper connections may have been forged through semantically connecting monolingual embeddings. Further analyses on this is provided in Appendix D.

6 Conclusion

Code-switching has gained research attention in the field of Natural Language Processing, especially for low-resourced languages due to most communities being largely multilingual. Large pre-trained multilingual models are typically a de-facto processing tool for many downstream tasks due to their increased processing capabilities. However, the use of explicit cross-lingual embeddings to bridge the gap between the language representations of the known codes lags behind. In this study, we explore the use of cross-lingual embeddings, that aim to bring known language pairs closer together to efficiently learn code-switching detection. Indeed, our experimental analyses show that mapping the switched languages into a single shared vector space before training shows un-tapped processing capabilities for code-switching detection. Concretely, fine-tuning AfroXLM-r, and Serengeti architectures with explicit cross-lingual embeddings outperforms monolingual (only joined through concatenation) embeddings. Although the cross-lingual injection experiments performed poorly compared to the original multilingual embedding baselines, our results imply that, learning explicit shared representation between known codes, formed a representational space enhancing inter-learning between token subspaces allowing efficient processing of code-switched text.

Transfer	Code-Switching dataset with base (eng) and second label switched (swtc)						
	xhozul	sottsn	xhozulsot	xhozultsn	sottsnxho	sottsnzul	xhozulsottsn
Afro-xlmr-b	86	93	95	96	94	94	95
Afro-xlmr-l	98	92	96	96	94	81	95
Serengeti	98	92	95	96	94	94	95

Table 4: Reports the label change transfer model’s F1-score for code-switching detection for each dataset averaged over 5 runs.

Models	Code-Switching dataset with base (eng)										
	xho	zul	sot	tsn	xhozul	sottsn	xhozulsot	xhozultsn	sottsnxho	sottsnzul	xhozulsottsn
Embedding Dimension: 50											
CCA	82	85	77	80	57	57	74	62	67	78	79
VecMap	76	84	85	84	81	65	83	76	67	63	45
Muse	76	80	63	68	70	70	73	74	59	62	70
Embedding Dimension: 100											
CCA	81	83	76	67	69	75	74	76	80	78	57
VecMap	82	84	78	85	69	74	83	46	76	68	62
Muse	67	81	76	74	72	71	74	64	70	72	70
Embedding Dimension: 200											
CCA	71	61	78	72	68	67	49	71	80	72	60
VecMap	71	84	80	76	73	75	79	70	78	59	78
Muse	76	80	77	75	71	71	74	75	72	72	72

Table 5: Reports the Afro-XLMr-base F1-score for code-switching detection for each dataset averaged over 5 runs using the three embedding techniques. The injected embeddings were randomly selected from the full embedding matrix to match the vocabulary size of the transformer.

Metric	Code-Switching dataset with base (eng)										
	xho	zul	sot	tsn	xhozul	sottsn	xhozulsot	xhozultsn	sottsnxho	sottsnzul	xhozulsottsn
Embedding Dimension: 50											
acc	93	81	91	91	80	82	87	83	67	76	79
prec	71	56	61	69	56	64	63	61	47	57	59
rec	68	53	62	65	56	62	64	61	40	54	61
f1	69	54	61	67	56	62	64	61	43	55	60
Embedding Dimension: 100											
acc	93	87	72	79	81	75	81	74	78	76	75
prec	63	58	57	61	58	55	55	52	59	59	55
rec	64	56	44	54	60	55	49	51	61	55	52
f1	63	56	49	57	59	55	52	51	60	56	52
Embedding Dimension: 200											
acc	92	81	76	91	72	82	65	65	55	81	75
prec	66	56	56	72	46	63	45	44	37	61	56
rec	64	52	33	66	44	62	41	40	28	63	54
f1	65	53	41	69	45	61	42	41	31	62	54

Table 6: Reports the Afro-XLMr-base Accuracy (acc), Precision (prec), Recall (rec), and F1-score (f1) for code-switching detection for each dataset averaged over 5 runs using FastText Monolingual embeddings and Glove embeddings for English words. The injected embeddings were randomly selected from the full embedding matrix to match the vocabulary size of the transformer.

607 Limitations

608 The vocabulary of the embeddings is larger than
609 the transformer embeddings. This means the selec-
610 tion of the appropriate words to add to the model
611 embedding becomes crucial.

612 We explored a theoretically ideal scenario for
613 generating cross-lingual embeddings and have not
614 explored set-ups such as how many bilingual lexi-
615 cons signals are sufficient for generating the best
616 shared representations.

617 In line with the above limitation, this study, did
618 not explore hyperparameter fine-tuning for injected
619 cross-lingual embeddings. We advise future works
620 to consider this as it may significantly improve
621 outcomes.

622 Our experiments did not consider the many
623 massively pre-trained multilingual models, such
624 as RemBert, mBERT, including afro-centric pre-
625 trained models such as AfroLM, AfriBerta, etc.
626 This is due to limitations in compute power as well
627 as time-constraints as recreating and injecting ex-
628 plicit cross-lingual embeddings takes more time,
629 even in the advent of multiprocessing capabilities
630 to speed up the process.

631 Ethics Statement

632 **Data and Models Disclaimer** The datasets and
633 models used in this study are collected from pub-
634 licly available resources with no potential harm,
635 threats, and risks to society.

636 License

637 This document and all its artifacts is licensed
638 under CC BY-SA 4.0. To view a copy of this
639 license, visit: [https://creativecommons.org/](https://creativecommons.org/licenses/by-sa/4.0/)
640 /licenses/by-sa/4.0/

641 References

642 Ife Adebara, AbdelRahim Elmadany, Muhammad
643 Abdul-Mageed, and Alcides Alcoba Inciarte. 2023.
644 **SERENGETI: Massively multilingual language**
645 **models for Africa.** In *Findings of the Association for*
646 *Computational Linguistics: ACL 2023*, pages 1498–
647 1537, Toronto, Canada. Association for Computa-
648 tional Linguistics.

649 David Ifeoluwa Adelani, Jade Abbott, Graham Neubig,
650 Daniel D’souza, Julia Kreutzer, Constantine Lignos,
651 Chester Palen-Michel, Happy Buzaaba, Shruti Rijh-
652 wani, Sebastian Ruder, et al. 2021. Masakhaner:
653 Named entity recognition for african languages.
654 *Transactions of the Association for Computational*
655 *Linguistics*, 9:1116–1131.

Gustavo Aguilar, Fahad AlGhamdi, Victor Soto, Mona Diab, Julia Hirschberg, and Thamar Solorio. 2019. Named entity recognition on code-switched data: Overview of the calcs 2018 shared task. <i>arXiv preprint arXiv:1906.04138</i> .	656
Jesujoba O. Alabi, David Ifeoluwa Adelani, Marius Mosbach, and Dietrich Klakow. 2022. Adapting pre-trained language models to African languages via multilingual adaptive fine-tuning. In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 4336–4349, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.	657
Mikel Artetxe, Gorka Labaka, and Eneko Agirre. 2018. A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings. <i>arXiv preprint arXiv:1805.06297</i> .	658
Mohammed Attia, Younes Samih, and Wolfgang Maier. 2018. Ghht at calcs 2018: Named entity recognition for dialectal arabic using neural networks. In <i>Proceedings of the Third Workshop on Computational Approaches to Linguistic Code-Switching</i> , pages 98–102.	659
Utsab Barman, Joachim Wagner, and Jennifer Foster. 2016. Part-of-speech tagging of code-mixed social media content: Pipeline, stacking and joint modelling. In <i>Proceedings of the second workshop on computational approaches to code switching</i> , pages 30–39.	660
Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. <i>Transactions of the association for computational linguistics</i> , 5:135–146.	661
Gokul Chittaranjan, Yogarshi Vyas, Kalika Bali, and Monojit Choudhury. 2014. Word-level language identification using crf: Code-switching shared task report of msr india system. In <i>Proceedings of the first workshop on computational approaches to code switching</i> , pages 73–79.	662
Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. <i>arXiv preprint arXiv:1911.02116</i> .	663
Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. <i>arXiv preprint arXiv:2207.04672</i> .	664
Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. <i>arXiv preprint arXiv:1810.04805</i> .	665
	666
	667
	668
	669
	670
	671
	672
	673
	674
	675
	676
	677
	678
	679
	680
	681
	682
	683
	684
	685
	686
	687
	688
	689
	690
	691
	692
	693
	694
	695
	696
	697
	698
	699
	700
	701
	702
	703
	704
	705
	706
	707
	708
	709

710	Cheikh M Bamba Dione, David Adelani, Peter Nabende, Jesujoba Alabi, Thapelo Sindane, Happy Buzaaba, Shamsuddeen Hassan Muhammad, Chris Chinene Emezue, Perez Ogayo, Anuoluwapo Aremu, et al. 2023. Masakhapos: Part-of-speech tagging for typologically diverse african languages. <i>arXiv preprint arXiv:2305.13989</i> .	764
711		765
712		766
713		767
714		768
715		769
716		770
717	Roald Eiselen and Martin J Puttkammer. 2014. Developing text resources for ten south african languages. In <i>LREC</i> , pages 3698–3703. Citeseer.	771
718		772
719		773
720	Manaal Faruqui and Chris Dyer. 2014. Improving vector space word representations using multilingual correlation. In <i>Proceedings of EACL</i> .	774
721		775
722		776
723	Aizaz Hussain and Muhammad Umair Arshad. 2021. An attention based neural network for code switching detection: English & roman urdu. <i>arXiv preprint arXiv:2103.02252</i> .	777
724		778
725		779
726		780
727	Navya Jose, Bharathi Raja Chakravarthi, Shardul Suryawanshi, Elizabeth Sherly, and John P McCrae. 2020. A survey of current datasets for code-switching research. In <i>2020 6th international conference on advanced computing and communication systems (ICACCS)</i> , pages 136–141. IEEE.	781
728		782
729		783
730		784
731		785
732		786
733	Amir Hossein Kargaran, Ayyoob Imani, François Yvon, and Hinrich Schütze. 2023. <i>Glotlid: Language identification for low-resource languages</i> . In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> .	787
734		788
735		789
736		790
737		791
738	Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. 2017. Unsupervised machine translation using monolingual corpora only. <i>arXiv preprint arXiv:1711.00043</i> .	792
739		793
740		794
741		795
742	Mack Makgatho, Vukosi Marivate, Tshephisho Sefara, and Valencia Wagner. 2021. Training cross-lingual embeddings for setswana and sepedi. <i>arXiv preprint arXiv:2111.06230</i> .	796
743		797
744		798
745		799
746	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. <i>Exploring the limits of transfer learning with a unified text-to-text transformer</i> . <i>arXiv e-prints</i> .	800
747		801
748		802
749		803
750	Thapelo Andrew Sindane and Vukosi Marivate. 2024. From n-grams to pre-trained multilingual models for language identification. In <i>Proceedings of the 4th International Conference on Natural Language Processing for Digital Humanities</i> , pages 229–239.	804
751		805
752		806
753	Kushagra Singh, Indira Sen, and Ponnurangam Kumaraguru. 2018. A twitter corpus for hindi-english code mixed pos tagging. In <i>Proceedings of the sixth international workshop on natural language processing for social media</i> , pages 12–17.	807
754		808
755		809
756	Sunayana Sitaram, Khyathi Raghavi Chandu, Sai Krishna Rallabandi, and Alan W Black. 2019. A survey of code-switched speech and language processing. <i>arXiv preprint arXiv:1904.00784</i> .	810
757		811
758		812
759		813
760		814
761	Thamar Solorio and Yang Liu. 2008. Part-of-speech tagging for english-spanish code-switched text. In <i>Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing</i> , pages 1051–1060.	815
762		816
763		817
		818
		819

Rob Van der Goot and Özlem Çetinoğlu.	2020.	Lexical normalization for code-switched data and its effect on pos-tagging.	<i>arXiv preprint arXiv:2006.01175</i> .	852
Pavanpankaj Vegi, Sivabhavani J, Biswajit Paul, Abhinav Mishra, Prashant Banjare, Prasanna Kumar K R, and Chitra Viswanathan.	2022.	Webcrawl african : A multilingual parallel corpora for african languages.	<i>In Proceedings of the Seventh Conference on Machine Translation</i> , pages 1076–1089, Abu Dhabi. Association for Computational Linguistics.	853
Yogarshi Vyas, Spandana Gella, Jatin Sharma, Kalika Bali, and Monojit Choudhury.	2014.	Pos tagging of english-hindi code-mixed social media content.	<i>In Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)</i> , pages 974–979.	854
Genta Indra Winata, Zhaojiang Lin, and Pascale Fung.	2019.	Learning multilingual meta-embeddings for code-switching named entity recognition.	<i>In Proceedings of the 4th Workshop on Representation Learning for NLP (RepL4NLP-2019)</i> , pages 181–186.	855
Genta Indra Winata, Chien-Sheng Wu, Andrea Madotto, and Pascale Fung.	2018.	Bilingual character representation for efficiently addressing out-of-vocabulary words in code-switching named entity recognition.	<i>arXiv preprint arXiv:1805.12061</i> .	856
Meng Xuan Xia.	2016.	Codeswitching language identification using subword information enriched word vectors.	<i>In Proceedings of the second workshop on computational approaches to code switching</i> , pages 132–136.	857
A Monolingual embeddings visualization				858
Visualizing the quality of monolingual embeddings can be a challenging task. In this Appendix, we present three visualization techniques: Principal Component Analyses (PCA), Uniform Manifold Approximation and Projection (UMAP), and t-distributed Stochastic Neighbor Embedding (t-SNE) for visualizing monolingual and cross-lingual embeddings in 2 dimensions. We present for each technique, visual plots for embeddings of 50, 100, 150, and 200 dimensions. For each pair of languages we concatenated the embeddings of the two languages and for words that appear in both vector spaces of the two languages we used the average of the two vectors as the vector representation of the duplicate word. Visualising all embeddings of word vectors in one plots causes difficulties in interpreting the word vectors relations and semantics, and such we created five clusters from the embeddings and only plotted a sample of 10 words from the clusters to get a sense of what the clusters contain. Figs. 3 to 14 , Figs. 16 to 27 , and Figs. 29 to 40 show the PCA, UMAP, and t-SNE plots when using monolingual embeddings respectively. Figs. 42 to 53 , Figs. 55 to 66 , and Figs. 68 to 79 show the PCA, UMAP, and t-SNE plots when using canonical correlation analysis (CCA) projection embeddings respectively, and Figs. 81 to 92 , Figs. 94 to 105 , and Figs. 107 to 118 show the PCA, UMAP, and t-SNE plots when using UMAP projection embeddings respectively. While Figs. 120 to 131 , Figs. 133 to 144 , and Figs. 146 to 157 show the PCA, UMAP, and t-SNE plots when using monolingual embeddings respectively. All three techniques were able to create clearly separable clusters with PCA showing more readable words within the clusters unlike UMAP and t-SNE. More importantly, word relations are captured within the embeddings clusters. For example, similar words such as ‘apere’ –translation ‘wore’, and ‘apara’ – translation ‘wear’; ‘babedi’ –translation ‘the two’ or ‘couple’, and ‘bedi’ translation ‘twice’, are captured in the same cluster.				859
B Cosine Similarities of Word Vectors				890
Measuring the similarity of word vectors for similar and dissimilar evaluation provides insights on the intrinsic quality of the embeddings. In this Appendix we used the available SimLex-999 dataset released by Makgatho et al. (2021) for measuring how closely related word vector are represented				891

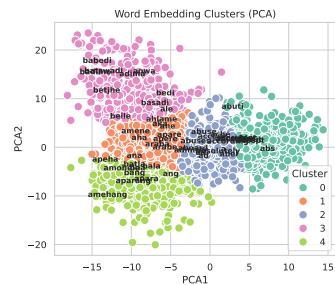


Figure 3: en-sot Emb

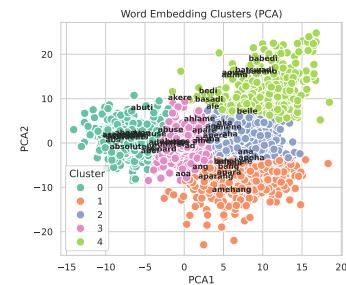


Figure 4: en-sot Emb

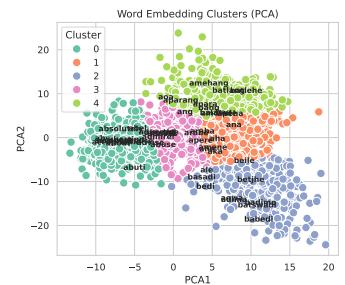


Figure 5: en-sot Emb

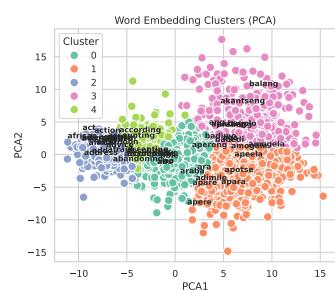


Figure 6: en-tsn Emb

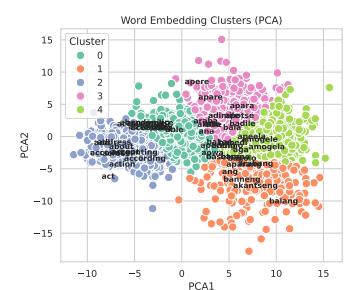


Figure 7: en-tsn Emb

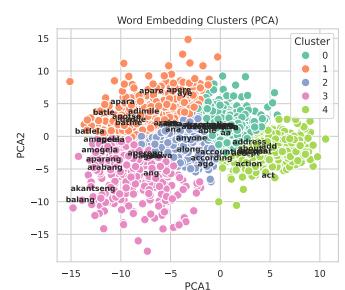


Figure 8: en-tsn Emb

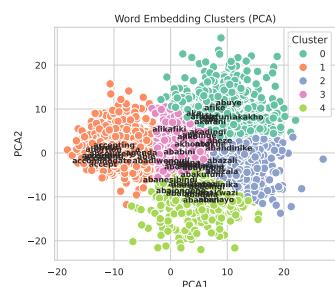


Figure 9: en-xho Emb

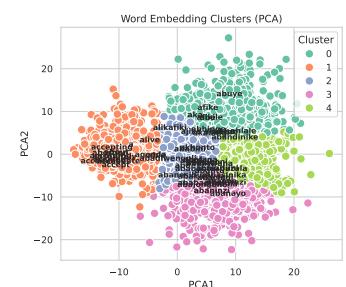


Figure 10: en-xho Emb

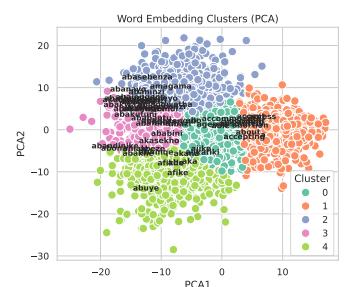


Figure 11: en-xho Emb

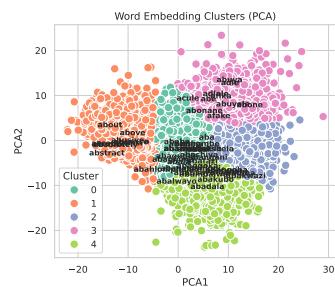


Figure 12: en-zul Emb

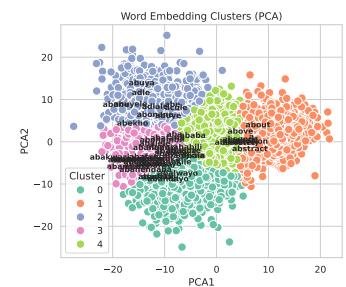


Figure 13: en-zul Emb

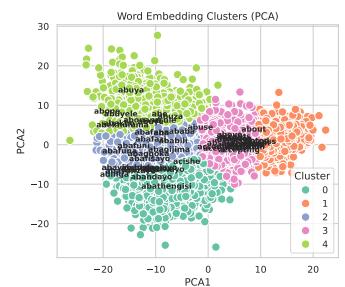


Figure 14: en-zul Emb

Figure 15: PCA Mono Emb plots for dimension 50 (left), 100(middle), and 200 (right)

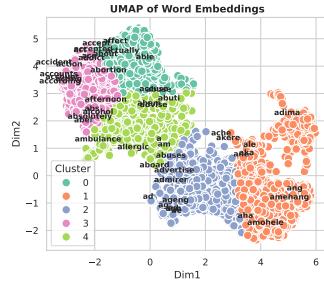


Figure 16: en-sot Emb

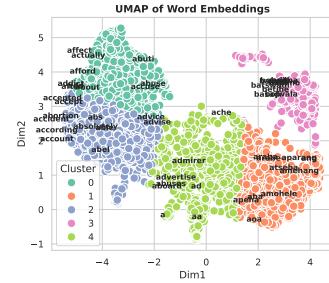


Figure 17: en-sot Emb

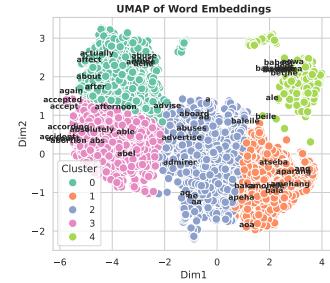


Figure 18: en-sot Emb

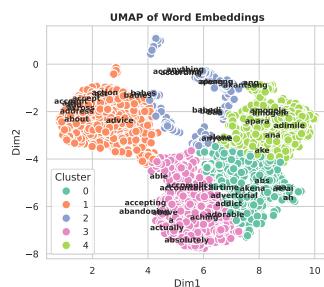


Figure 19: en-tsn Emb

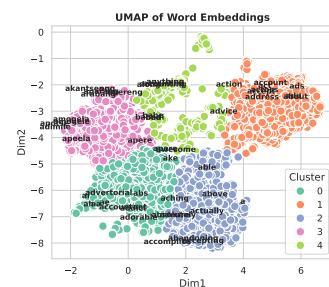


Figure 20: en-tsn Emb

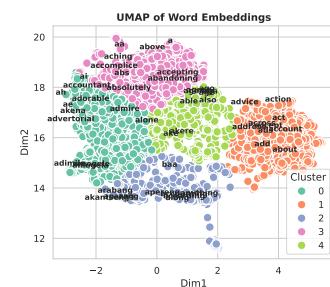


Figure 21: en-tsn Emb

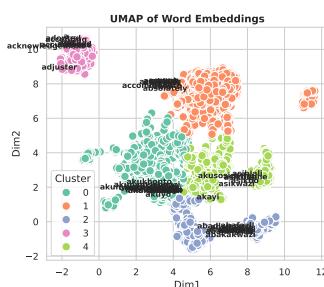


Figure 22: en-xho Emb

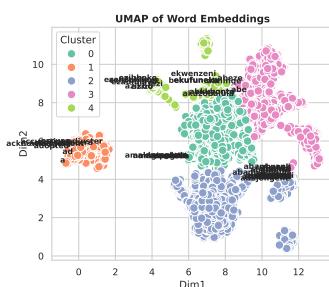


Figure 23: en-xho Emb

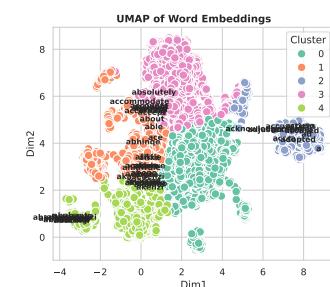


Figure 24: en-xho Emb

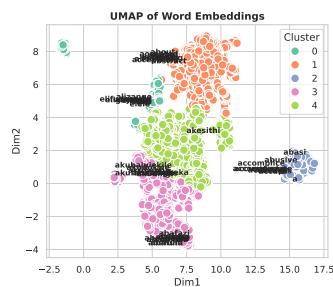


Figure 25: en-zul Emb

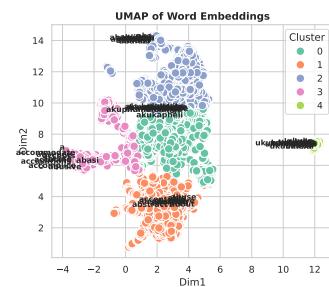


Figure 26: en-zul Emb

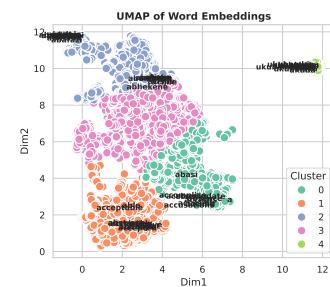


Figure 27: en-zul Emb

Figure 28: UMAP Mono Emb plots for dimension 50 (left), 100(middle), and 200 (right)

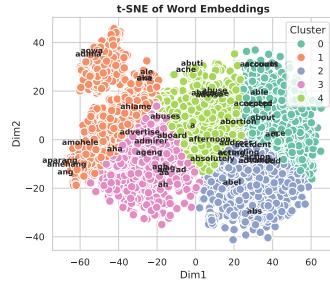


Figure 29: en-sot Emb

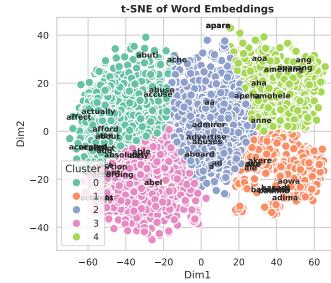


Figure 30: en-sot Emb

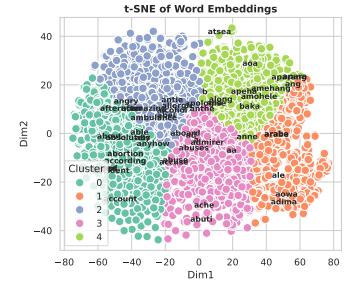


Figure 31: en-sot Emb

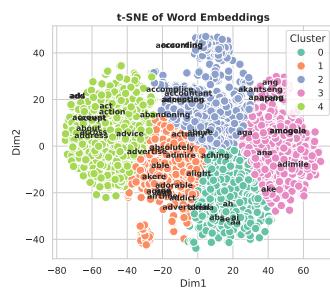


Figure 32: en-tsn Emb

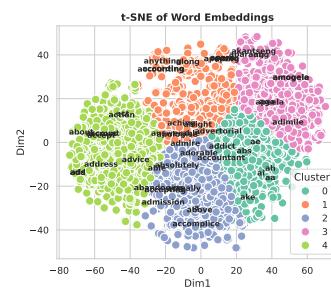


Figure 33: en-tsn Emb

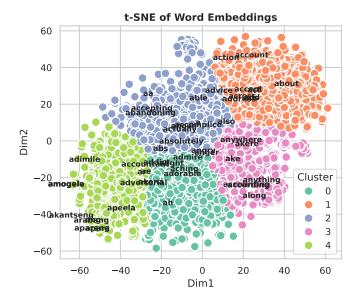


Figure 34: en-tsn Emb

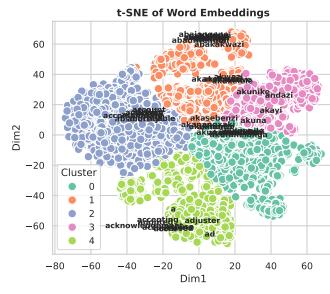


Figure 35: en-xho Emb

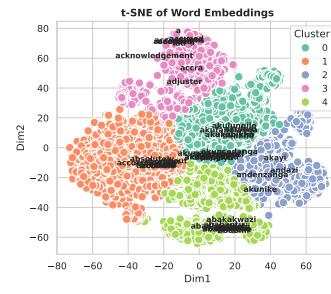


Figure 36: en-xho Emb

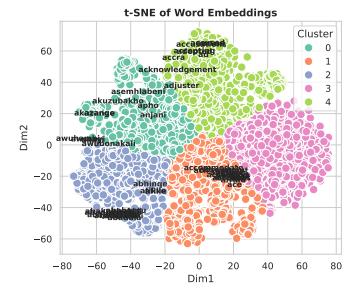


Figure 37: en-xho Emb

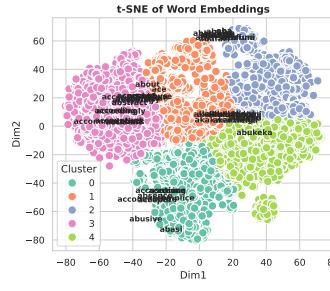


Figure 38: en-zul Emb

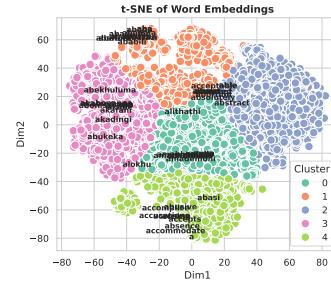


Figure 39: en-zul Emb

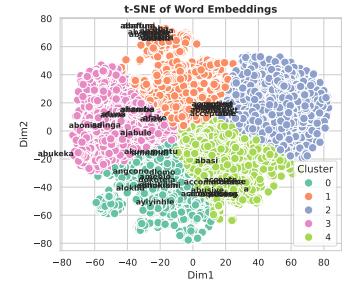


Figure 40: en-zul Emb

Figure 41: t-SNE Mono Emb plots for dimension 50 (left), 100(middle), and 200 (right)

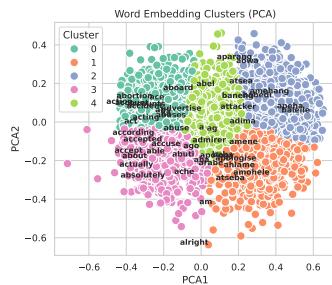


Figure 42: en-sot Emb

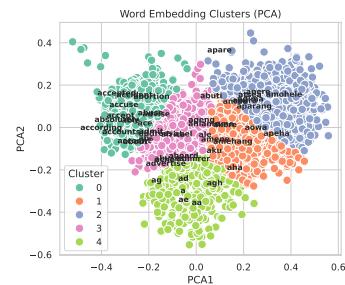


Figure 43: en-sot Emb

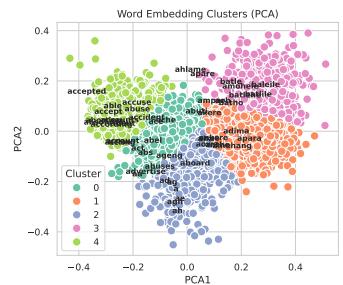


Figure 44: en-sot Emb

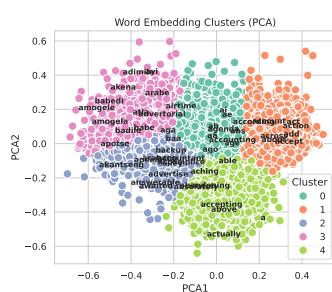


Figure 45: en-tsn Emb

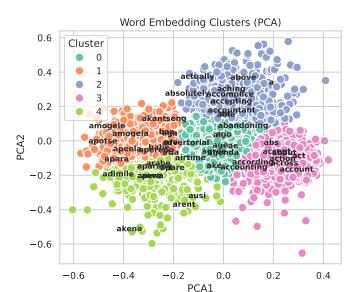


Figure 46: en-tsn Emb

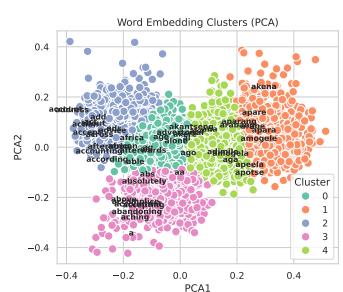


Figure 47: en-tsn Emb

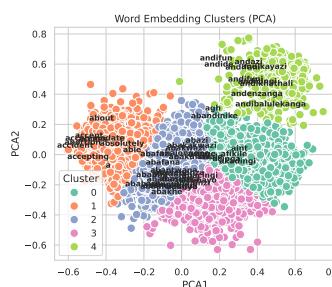


Figure 48: en-xho Emb

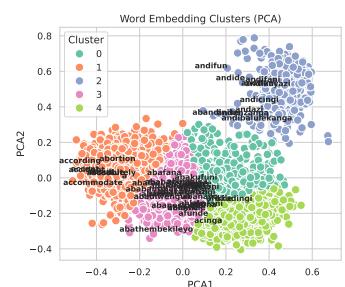


Figure 49: en-xho Emb

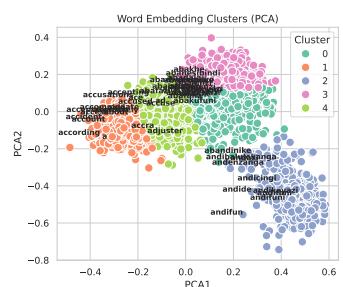


Figure 50: en-xho Emb

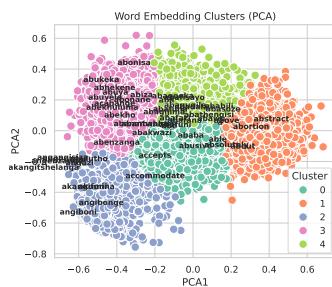


Figure 51: en-zul Emb

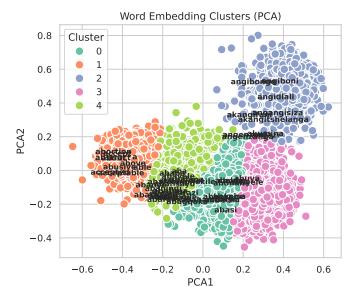


Figure 52: en-zul Emb

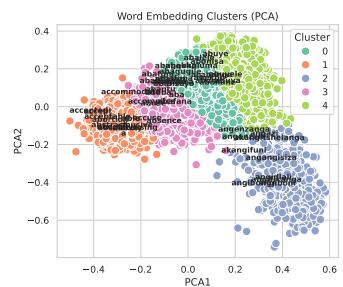


Figure 53: en-zul Emb

Figure 54: PCA CCA Emb plots for dimension 50 (left), 100(middle), and 200 (right)

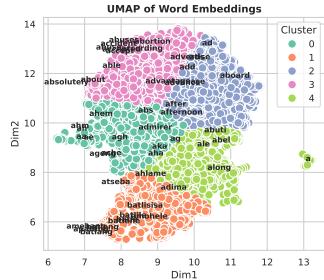


Figure 55: en-sot Emb

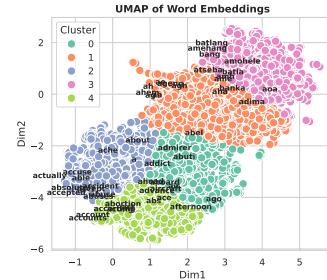


Figure 56: en-sot Emb

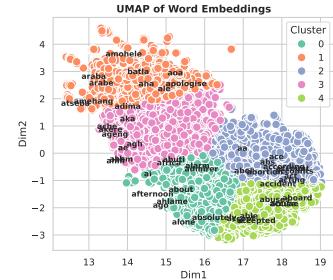


Figure 57: en-sot Emb

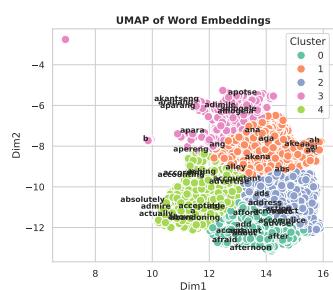


Figure 58: en-tsn Emb

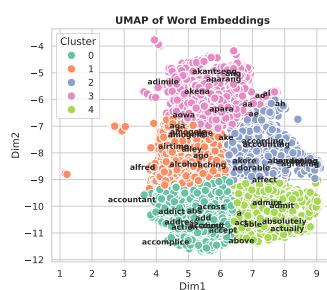


Figure 59: en-tsn Emb

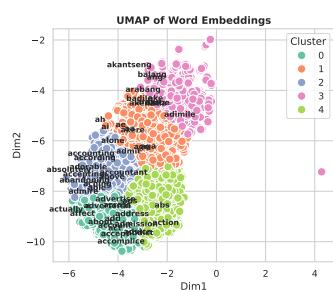


Figure 60: en-tsn Emb

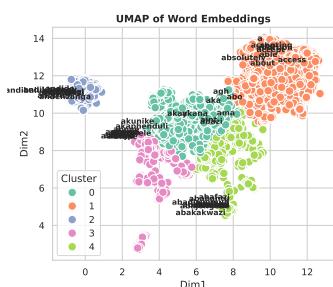


Figure 61: en-xho Emb

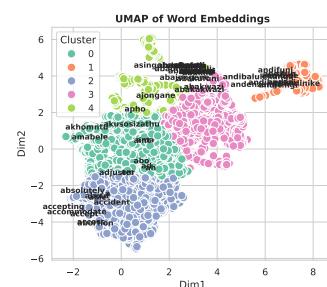


Figure 62: en-xho Emb

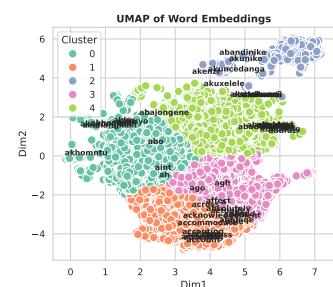


Figure 63: en-xho Emb

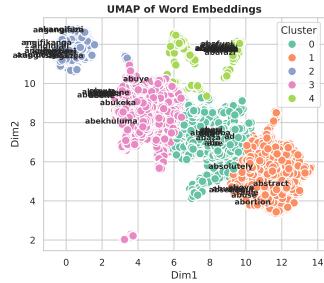


Figure 64: en-zul Emb

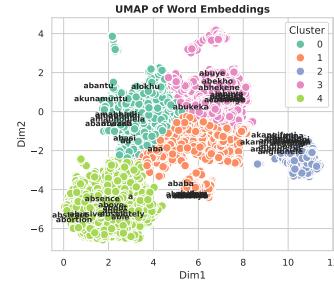


Figure 65: en-zul Emb

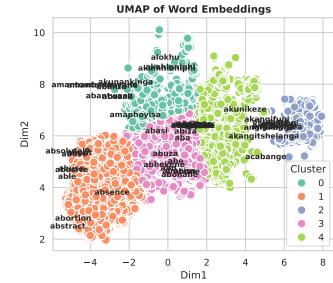


Figure 66: en-zul Emb

Figure 67: UMAP CCA Emb plots for dimension 50 (left), 100(middle), and 200 (right)

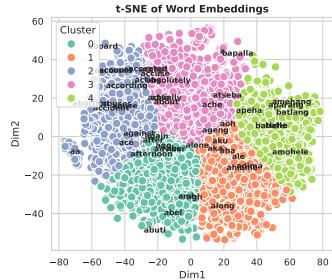


Figure 68: en-sot Emb

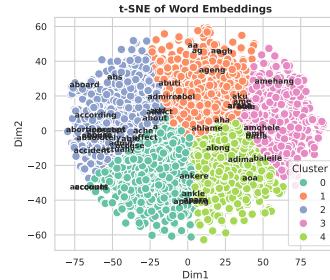


Figure 69: en-sot Emb

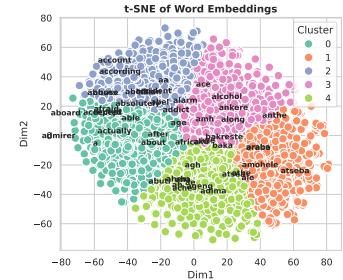


Figure 70: en-sot Emb

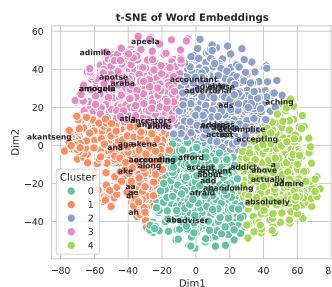


Figure 71: en-tsn Emb

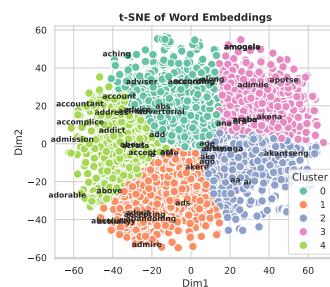


Figure 72: en-tsn Emb

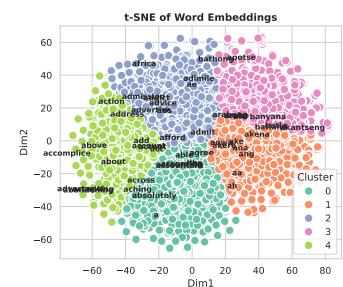


Figure 73: en-tsn Emb

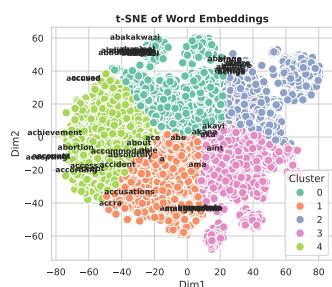


Figure 74: en-xho Emb

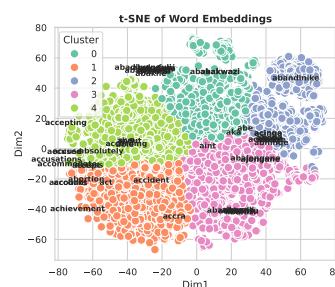


Figure 75: en-xho Emb

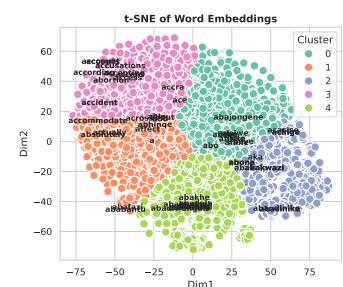


Figure 76: en-xho Emb

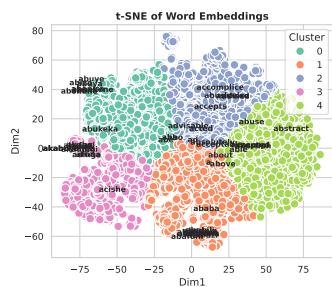


Figure 77: en-zul Emb

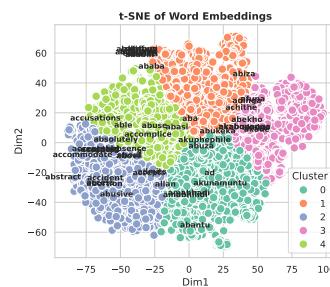


Figure 78: en-zul Emb

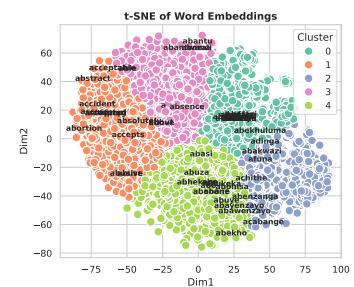


Figure 79: en-zul Emb

Figure 80: t-SNE CCA Emb plots for dimension 50 (left), 100(middle), and 200 (right)

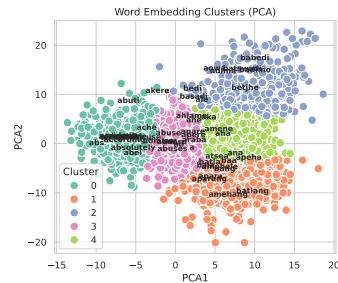


Figure 81: en-sot Emb

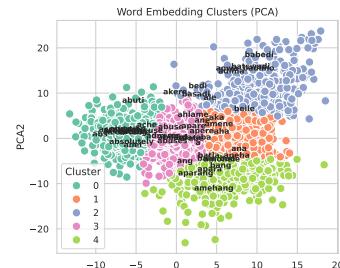


Figure 82: en-sot Emb

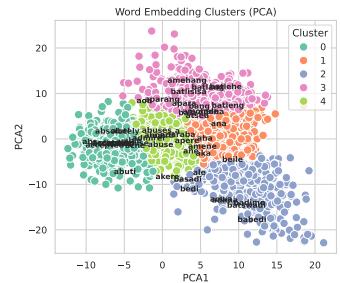


Figure 83: en-sot Emb

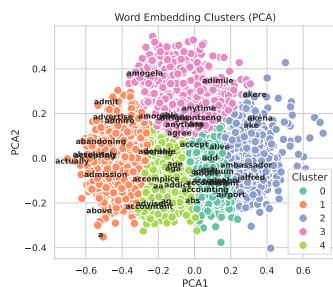


Figure 84: en-xho Emb

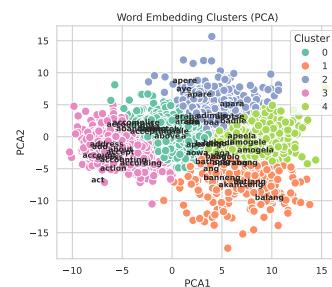


Figure 85: en-xho Emb

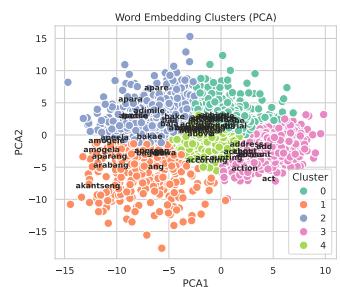


Figure 86: en-tsn Emb

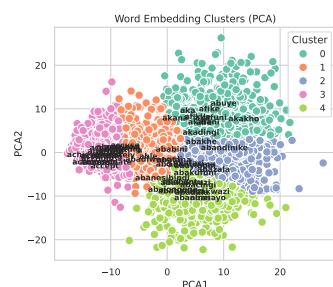


Figure 87: en-xho Emb

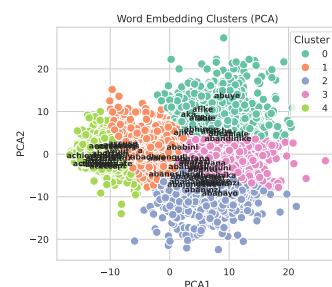


Figure 88: en-xho Emb

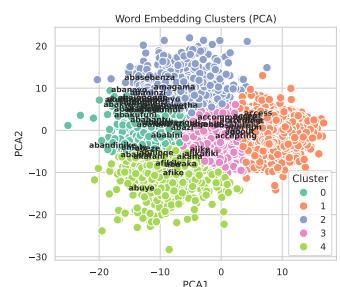


Figure 89: en-xho Emb

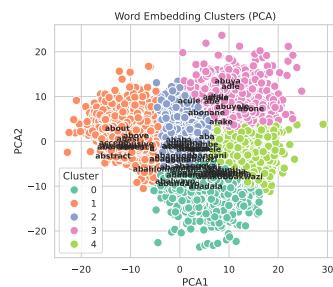


Figure 90: en-zul Emb

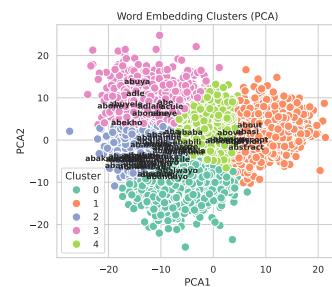


Figure 91: en-zul Emb

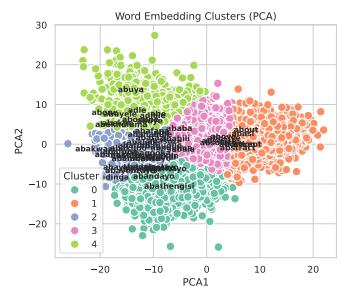


Figure 92: en-zul Emb

Figure 93: PCA Muse Emb plots for dimension 50 (left), 100(middle), and 200 (right)

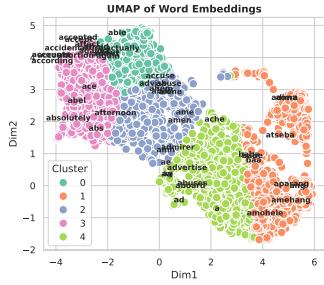


Figure 94: en-sot Emb

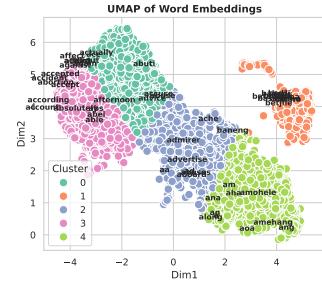


Figure 95: en-sot Emb

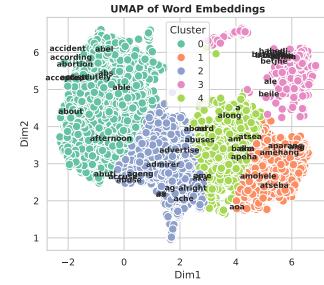


Figure 96: en-sot Emb

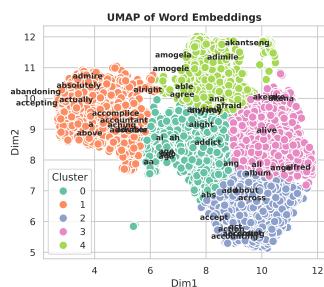


Figure 97: en-tsn Emb

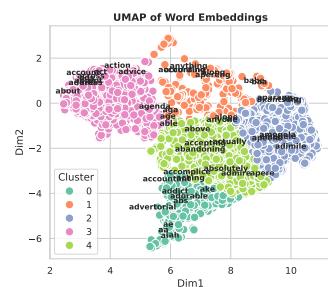


Figure 98: en-tsn Emb

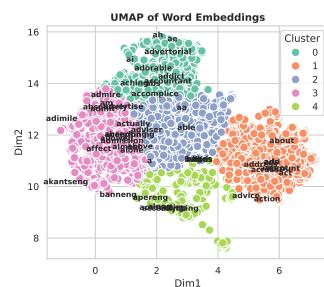


Figure 99: en-tsn Emb

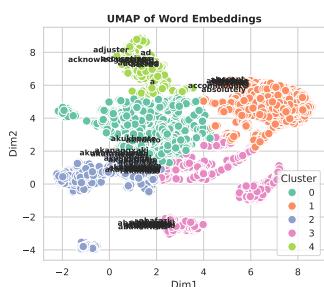


Figure 100: en-xho Emb

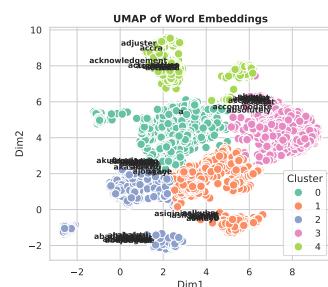


Figure 101: en-xho Emb

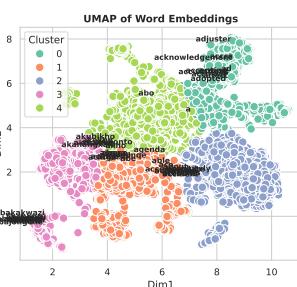


Figure 102: en-xho Emb

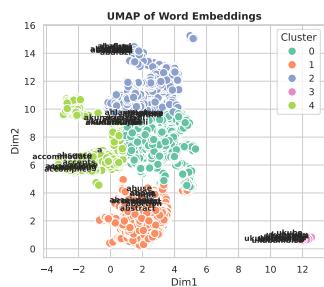


Figure 103: en-zul Emb

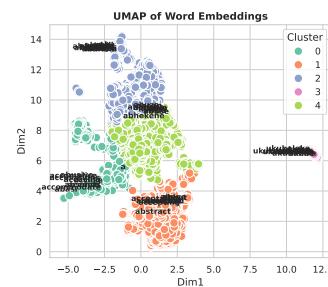


Figure 104: en-zul Emb

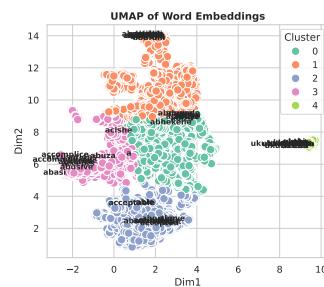


Figure 105: en-zul Emb

Figure 106: UMAP Muse Emb plots for dimension 50 (left), 100(middle), and 200 (right)

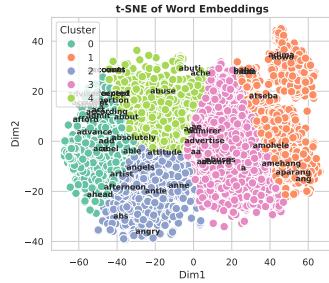


Figure 107: en-sot Emb

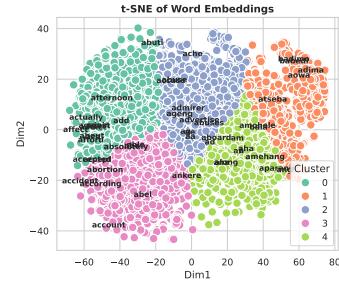


Figure 108: en-sot Emb

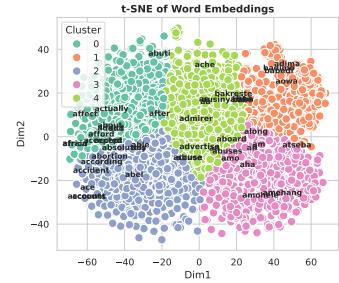


Figure 109: en-sot Emb

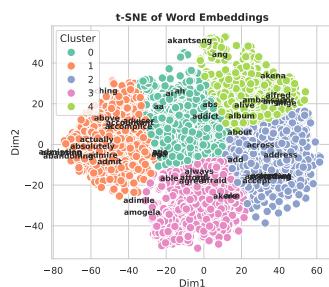


Figure 110: en-tsn Emb

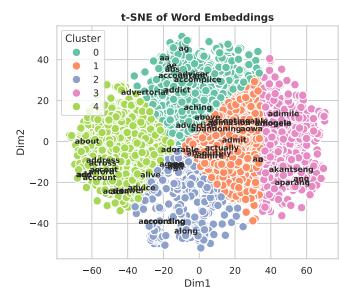


Figure 111: en-tsn Emb

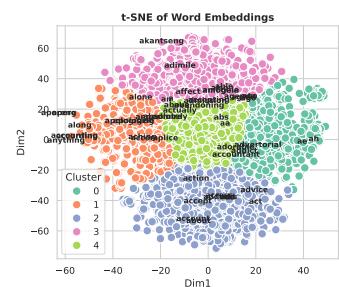


Figure 112: en-tsn Emb

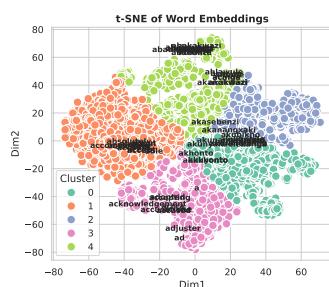


Figure 113: en-zul Emb

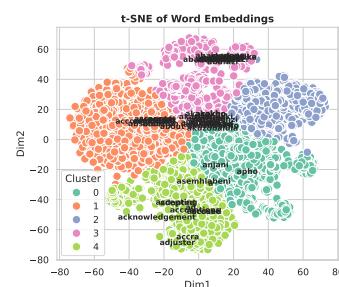


Figure 114: en-zul Emb

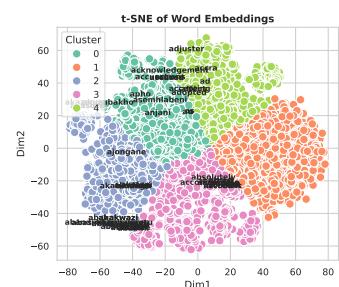


Figure 115: en-zul Emb

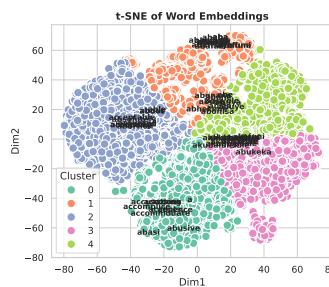


Figure 116: en-zul Emb

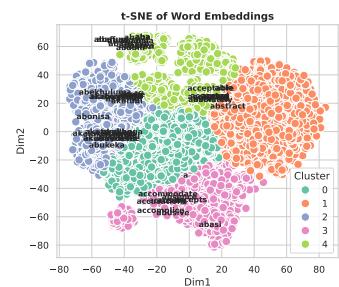


Figure 117: en-zul Emb

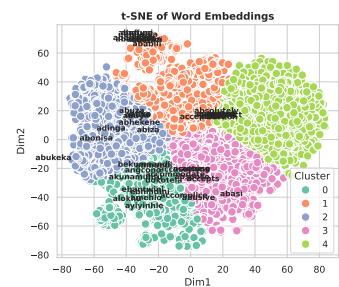


Figure 118: en-zul Emb

Figure 119: t-SNE Muse Emb plots for dimension 50 (left), 100(middle), and 200 (right)

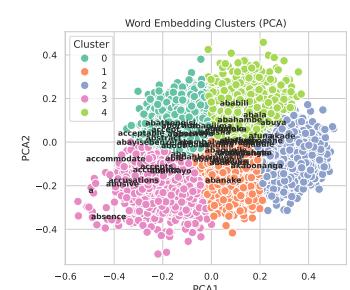
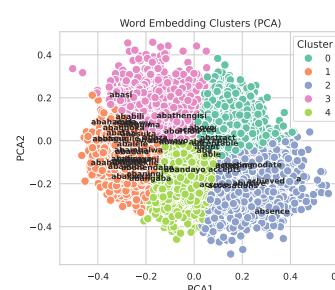
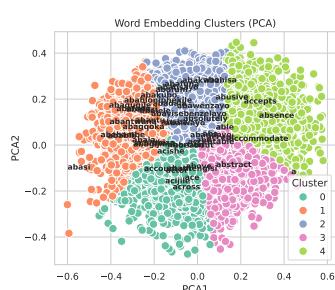
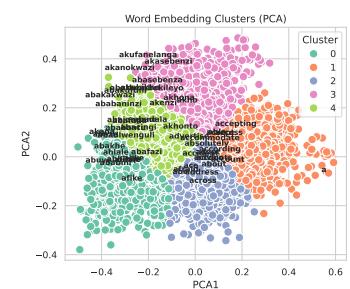
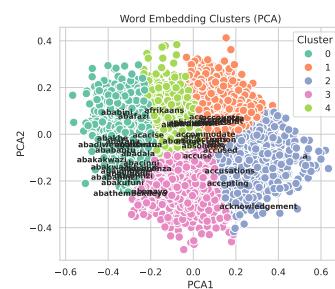
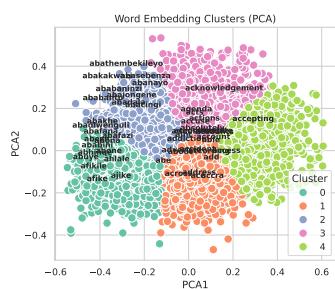
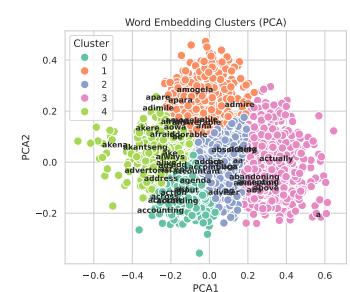
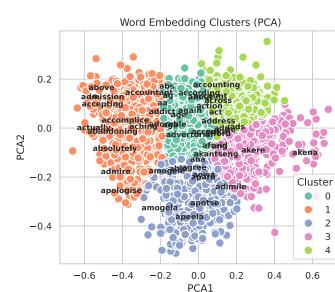
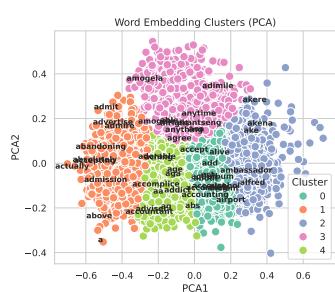
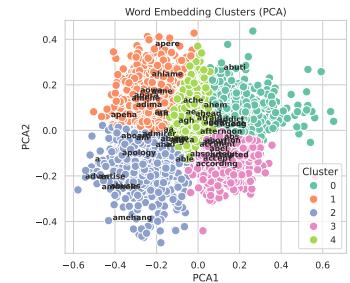
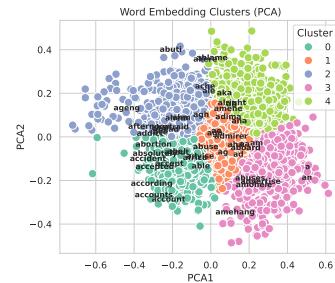
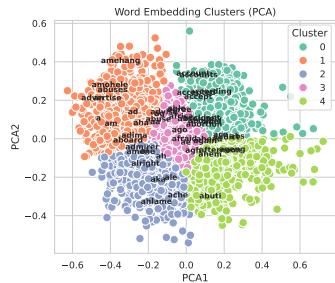


Figure 132: PCA VecMap Emb plots for dimension 50 (left), 100(middle), and 200 (right)

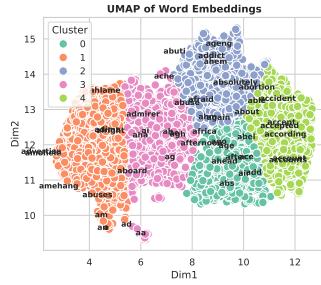


Figure 133: en-sot Emb

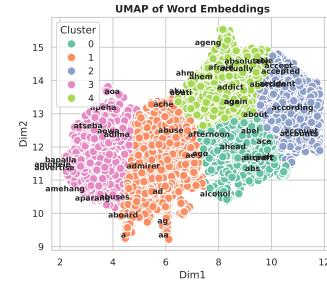


Figure 134: en-sot Emb

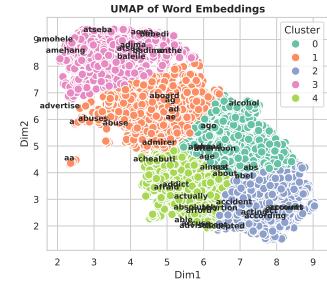


Figure 135: en-sot Emb

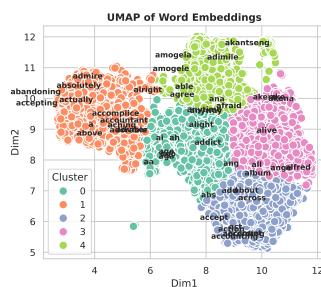


Figure 136: en-tsn Emb

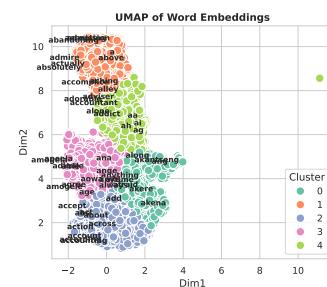


Figure 137: en-tsn Emb

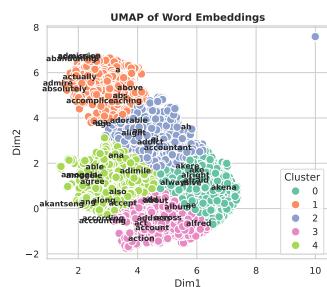


Figure 138: en-tsn Emb

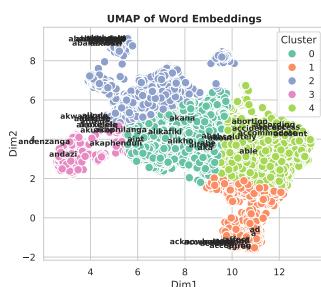


Figure 139: en-xho Emb

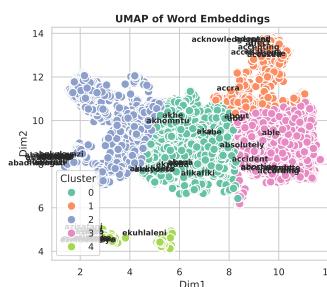


Figure 140: en-xho Emb

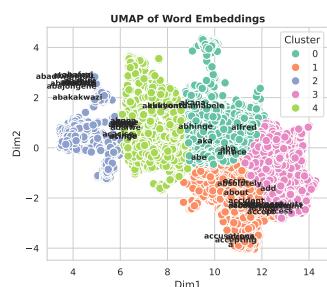


Figure 141: en-xho Emb

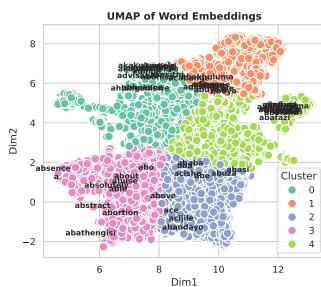


Figure 142: en-zul Emb

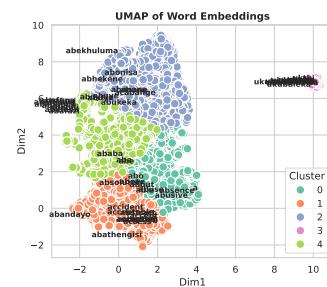


Figure 143: en-zul Emb

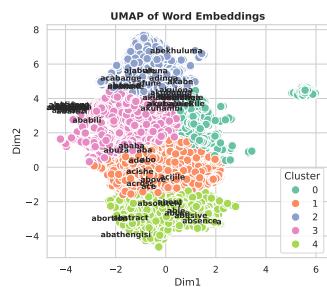


Figure 144: en-zul Emb

Figure 145: UMAP VecMap Emb plots for dimension 50 (left), 100(middle), and 200 (right)

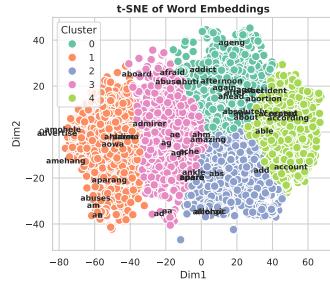


Figure 146: en-sot Emb

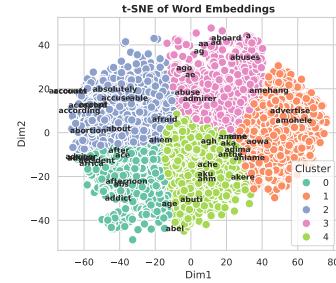


Figure 147: en-sot Emb

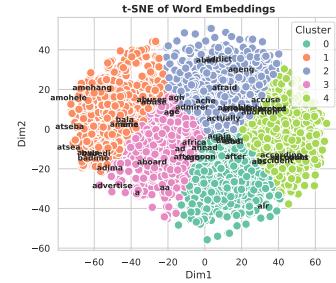


Figure 148: en-sot Emb

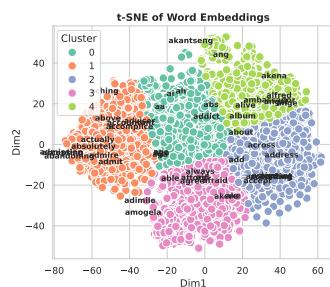


Figure 149: en-tsn Emb

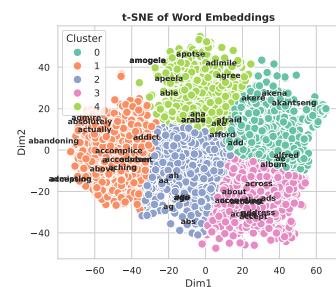


Figure 150: en-tsn Emb

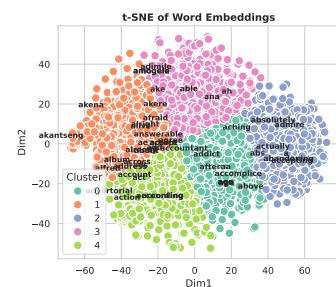


Figure 151: en-tsn Emb

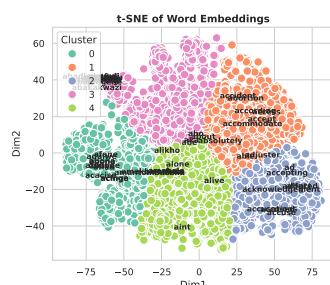


Figure 152: en-xho Emb

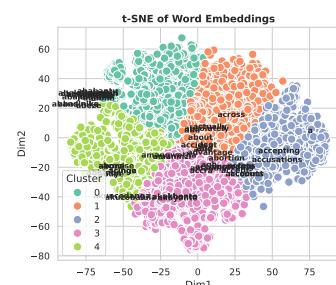


Figure 153: en-xho Emb

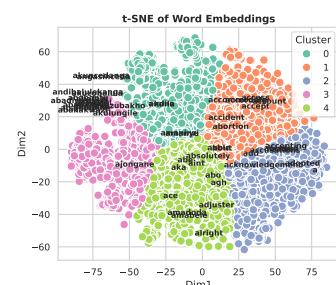


Figure 154: en-xho Emb

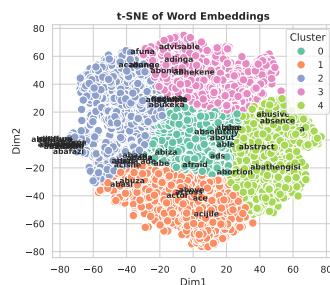


Figure 155: en-zul Emb

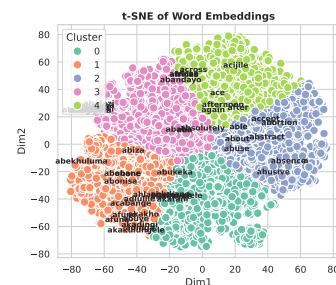


Figure 156: en-zul Emb

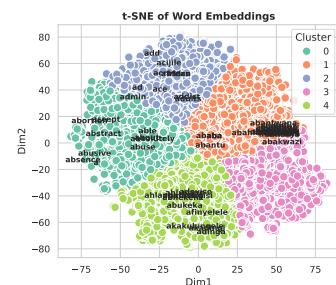


Figure 157: en-zul Emb

Figure 158: t-SNE VecMap Emb plots for dimension 50 (left), 100(middle), and 200 (right)

in word embedding spaces. The dataset has word pairs for Sepedi and Setswana only. The datasets is derived from the original English SimLex-999 in which the original English words (word1 and word2) are translated to translate1 and translate2 for the two languages Sesotho and Setswana. That is, each dataset will contain four pairs - 2 original English pairs and the 2 translate pairs together with a score for each pair to measure the similarity of the words. For our experimental purposes, we investigated the similarity of the combinations – low-resourced to low-resourced for the translated words (lr-t), high-resourced to high-resourced (hr), high-resourced to low-resourced translate 2 (hrlr-t), and low-resourced translate 1 to high-resourced word 2 (lr-rhr). Our monolingual embeddings of the four languages Setswana, Sesotho, IsiXhosa and Isizulu were evaluated on the two available datasets of Sesotho, and Setswana. The following Section B.1, and B.2 reports the Cosine plots for the aforementioned datasets.

B.1 Monolingual cosine similarity scores

Figs. 159 to 183, Figs. 185 to 209, Figs. 211 to 235, and Figs. 237 to 261 show the monolingual cosine similarity plots for four languages tsn, sot, xho, and zul, including their different dimension (50, 100, and 200). For each dataset lr-t, hr, hrlr-t, and lr-thr, we sampled 10 paired words using random sample 10, and generated the cosine plots for the 10 pairs. These plots rigorously show across multiple setups that monolingual embeddings where not able o capture word similarities sufficiently, especially when considering similar words from different languages. Cross-lingual embeddings, improves on this limitations and this is illustrated in the next subsection.

B.2 Cross-lingual cosine similarity scores

Similar to Section B.1, for each dataset lr-t, hr, hrlr-t, and lr-thr, 10 samples were extracted and the cosine similarity scores for the vector representations of the pairs were calculated. To analyse the intrinsic quality of monolingual embeddings, Figs. 159 to 261 shows the cosine similarity of paired words calculated from their vectors representations extracted from the monolingual embeddings. These figures not only show that the monolingual embeddings for the considered low-resourced languages are not able to capture similar words within a single language (i.e monolingual similarity) effectively but also fail to capture similarities across languages (cross-lingual similarity). This is could be the im-

pact of insufficient training data for low-resourced languages, which is not the case for English Glove embeddings as reflected by high monolingual similarity scores. However, using cross-lingual embeddings derived using the three projection techniques: Canonical Correlation Analyses (CCA), VecMap, and Muse, improvements of the two similarities modes were observed. Figs. 263 to 365, represent cosine similarity plots for CCA, Figs. 367 to 469, show cosine plots for Muse embeddings, while Figs. 471 to 573, show plots for Muse projection techniques. Interestingly, closely similar word such as 'bonolo' in Setswana - translation 'easy', and the word 'simple', were captured with high scores while similarities scores of words in the monolingual case appear to not capture any meaning (e.g sometimes very high for very distant words). This means that, in the cross-lingual case, the independently trained vector spaces are somehow projected into a common shared space where geometrical distance of the vectors captures word similarity and meaning even for words not belonging to the same language.

B.3 Spearman's correlations

The SimLex-999 dataset, unlike other word similarity datasets, measures word similarity outside of relatedness or association. This is often measured using Spearman's correlation of the cosine similarity scores of the word pairs and the human-annotated scores. In this Appendix, we measured the Spearman's correlation of the cosine scores for our Setswana (tsn), and Sesotho (sot) monolingual embeddings (see Figs. 575 to 587 respectively). We only considered tsn, and sot embeddings since these are the two languages closest to our available test datasets of tsn and nso released by Makgatho et al. (2021). The Spearman scores generated using monolingual embeddings indicate low values for all datasets and all embedding dimensions. This is expected since the cosine scores are generated from independently trained English and the paired low-resourced language embeddings. Moreover, we measured the Spearman's correlation using cross-lingual embeddings generated with CCA Figs. 589 to 601, Muse Figs. 603 to 615, and VecMap Figs. 617 to 629. Unlike heatmaps in the monolingual case, cross-lingual embeddings generate spread-out heatmaps with increased scores across languages. Nevertheless, significant improvement is realized on the Sesotho (sot), and

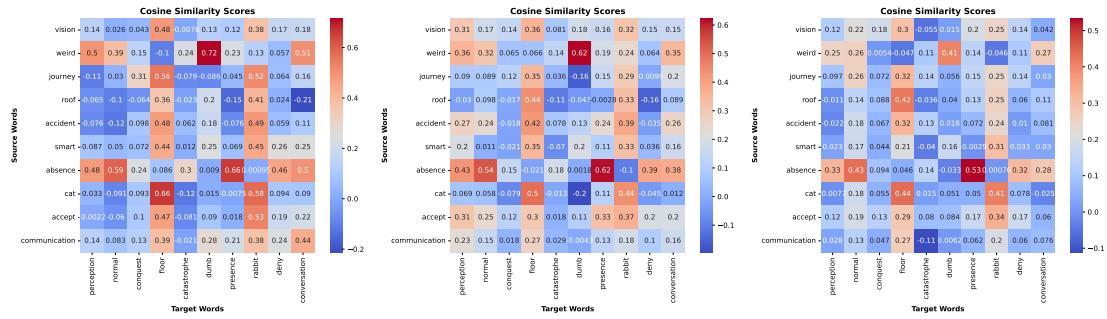


Figure 159: hr

Figure 160: hr

Figure 161: hr

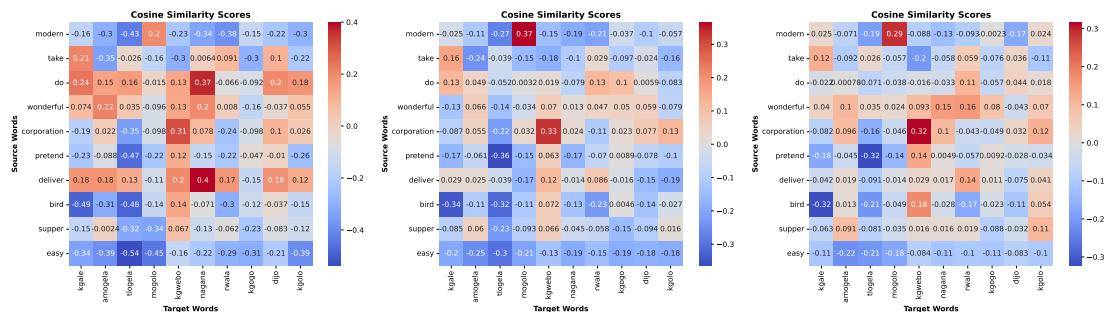


Figure 162: hr and lr-t

Figure 163: hr and lr-t

Figure 164: hr and lr-t

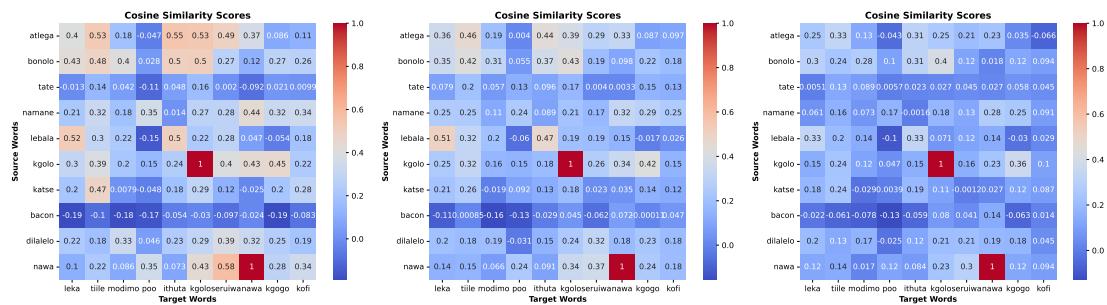


Figure 165: lr-t

Figure 166: lr-t

Figure 167: lr-t

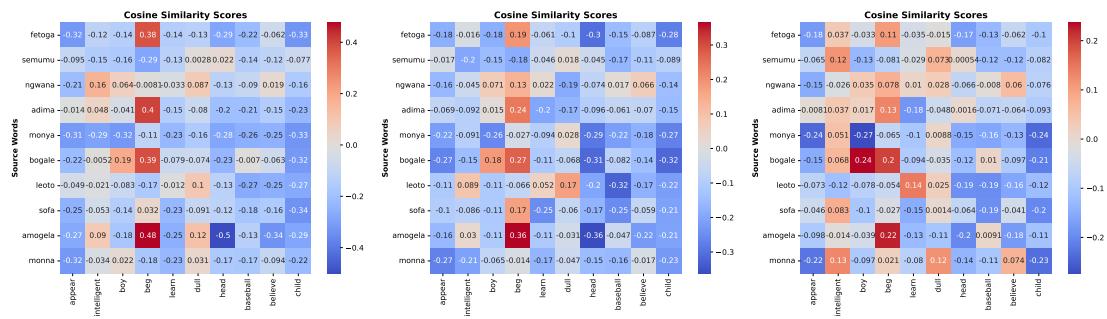


Figure 168: lr-t and hr

Figure 169: lr-t and hr

Figure 170: lr-t and hr

Figure 171: sot Mono Emb plots of 50 (left), 100(middle), and 200 (right), nso test

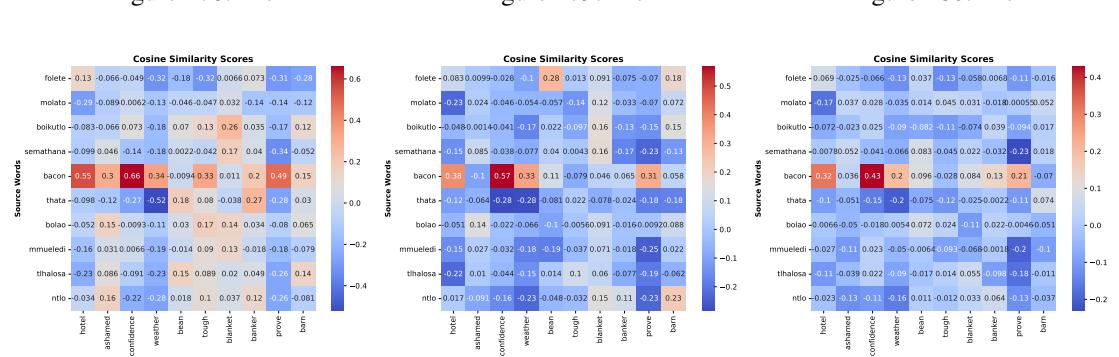
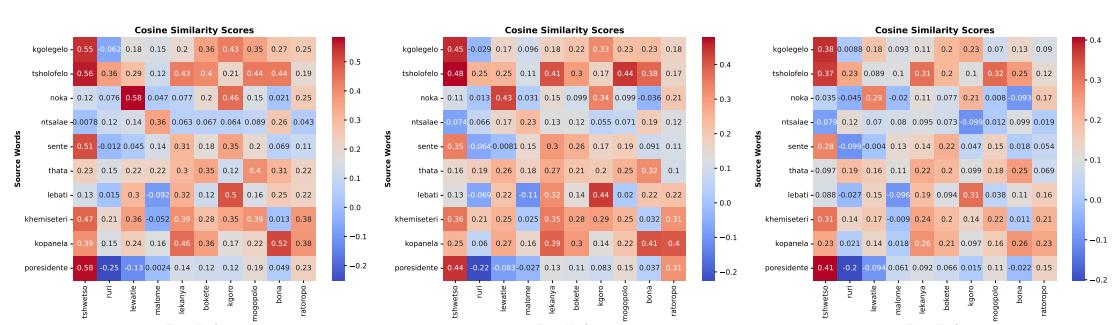
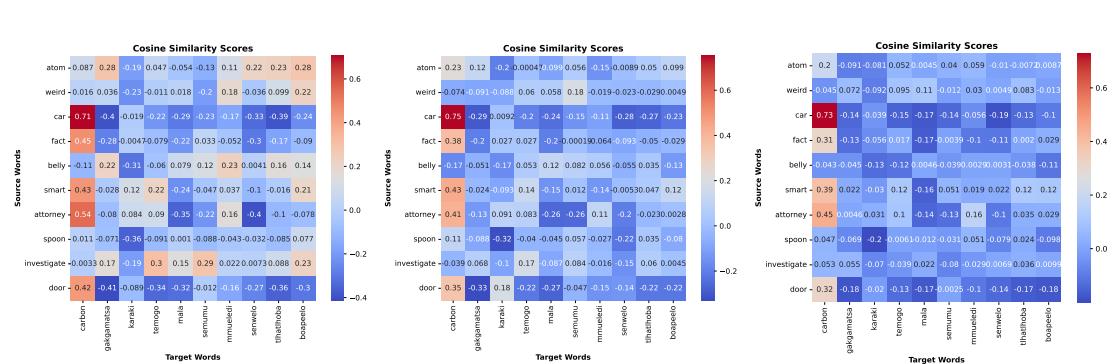
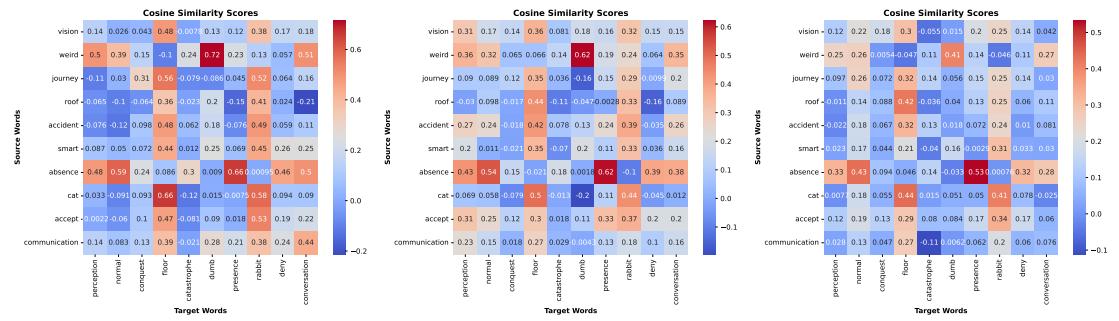


Figure 184: tsn Mono Emb plots of 50 (left), 100(middle), and 200 (right), tsn test

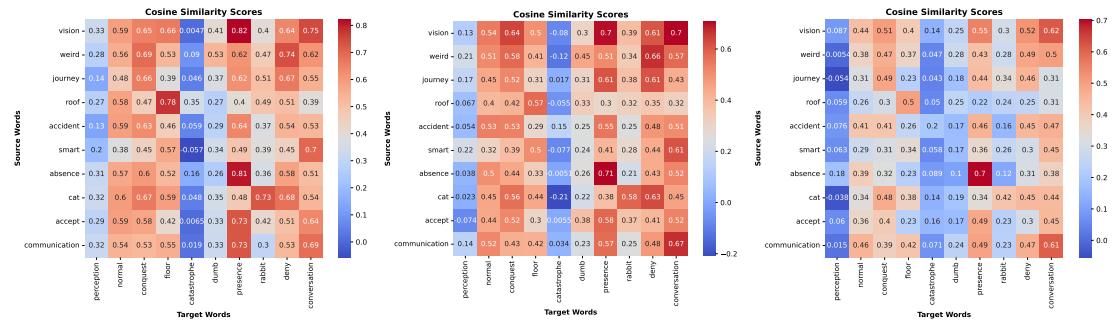


Figure 185: hr

Figure 186: hr

Figure 187: hr

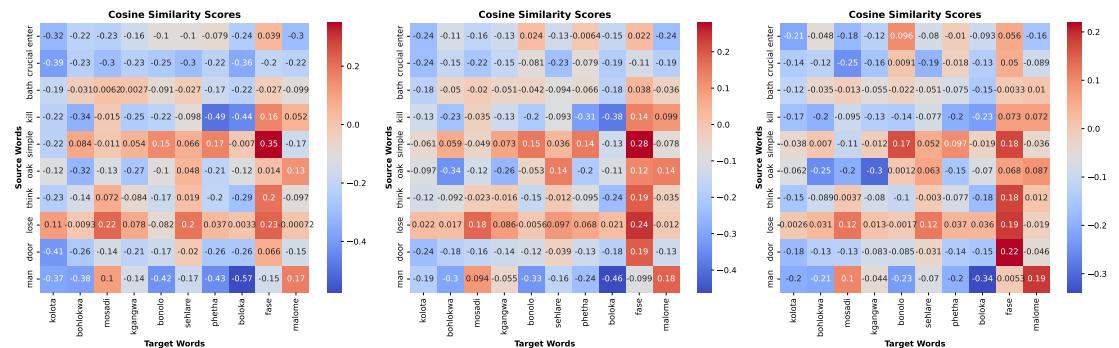


Figure 188: hr and lr-t

Figure 189: hr and lr-t

Figure 190: hr and lr-t

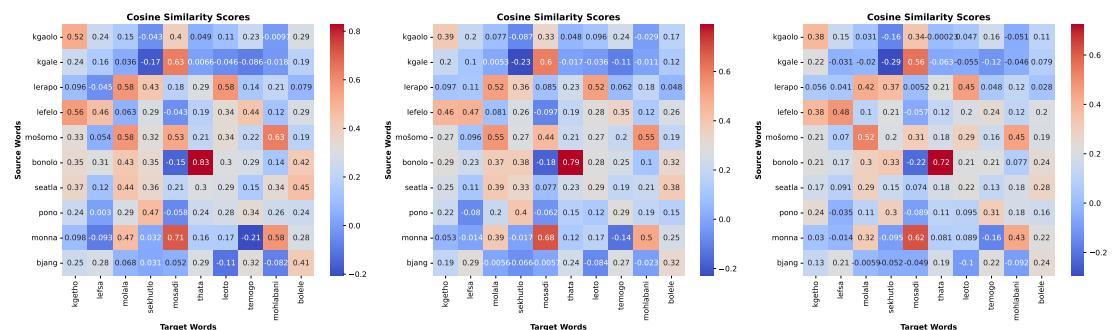


Figure 191: lr-t

Figure 192: lr-t

Figure 193: lr-t

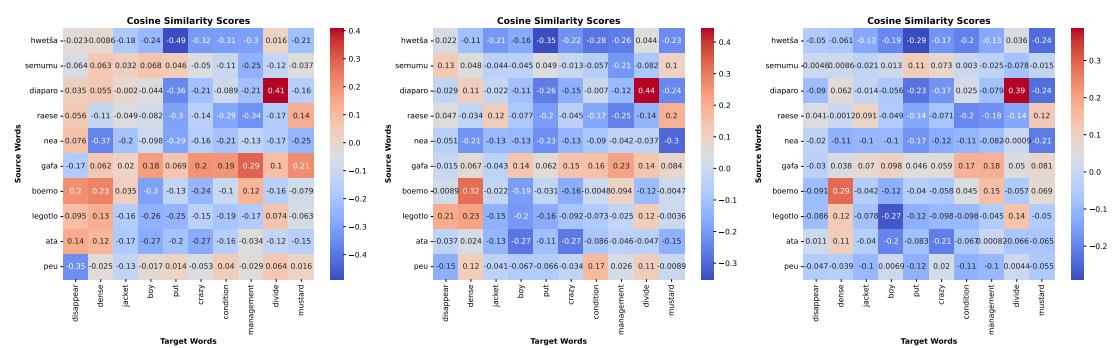


Figure 194: lr-t and hr

Figure 195: lr-t and hr

Figure 196: lr-t and hr

Figure 197: sot Mono Emb plots of 50 (left), 100(middle), and 200 (right), nso test

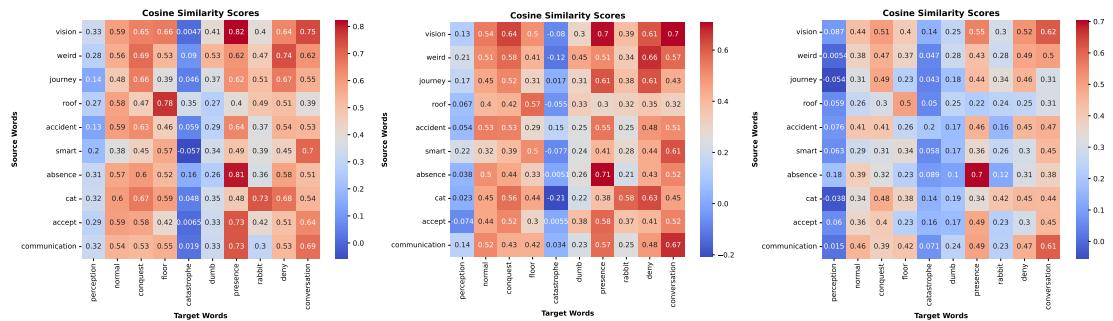


Figure 198: hr

Figure 199: hr

Figure 200: hr

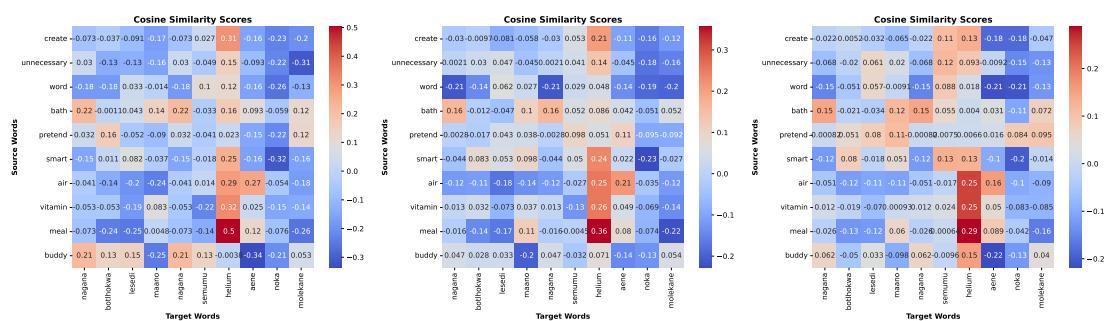


Figure 201: hr and lr-t

Figure 202: hr and lr-t

Figure 203: hr and lr-t

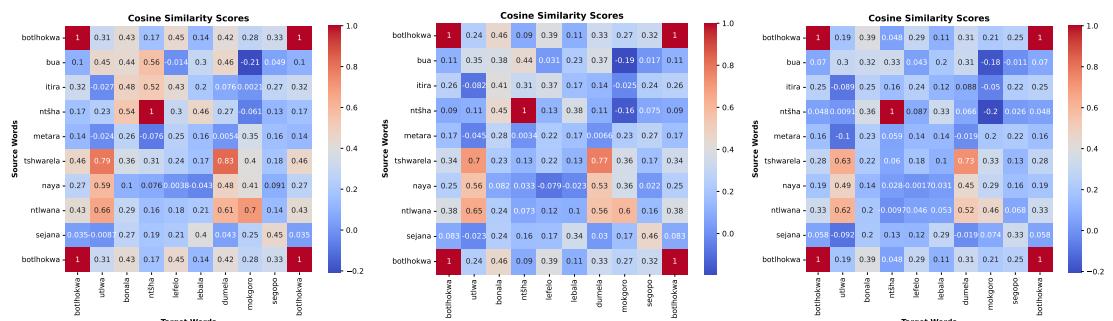


Figure 204: lr-t

Figure 205: lr-t

Figure 206: lr-t

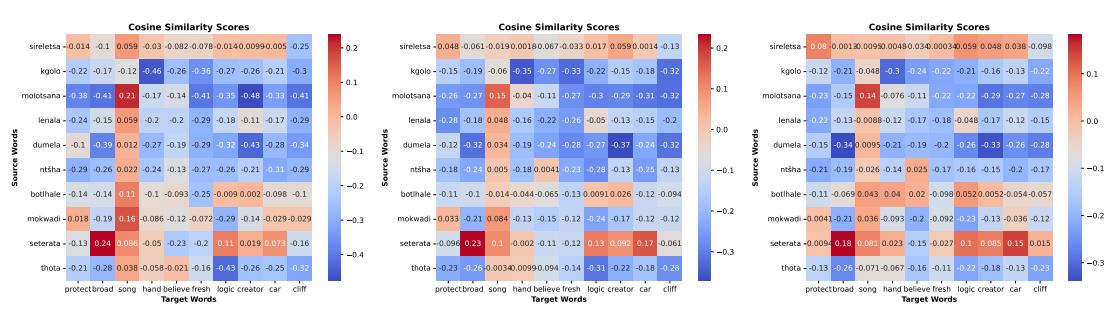


Figure 207: lr-r and hr

Figure 208: lr-t and hr

Figure 209: lr-t and hr

Figure 210: sot Mono Emb plots of 50 (left), 100(middle), and 200 (right), tsn test

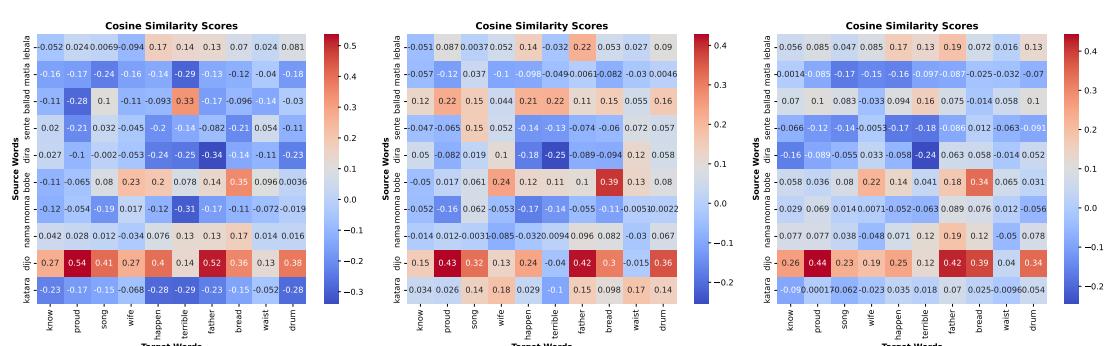
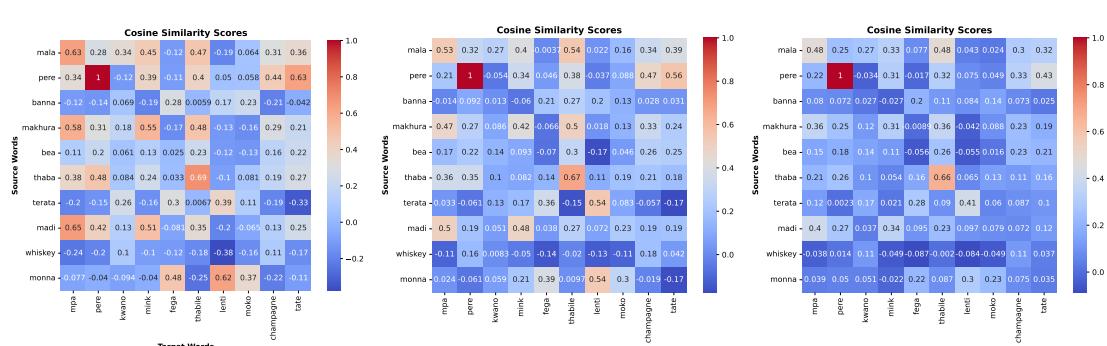
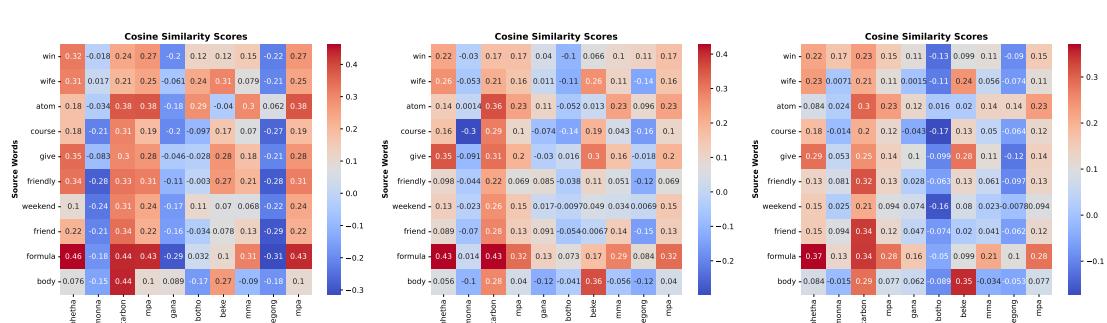
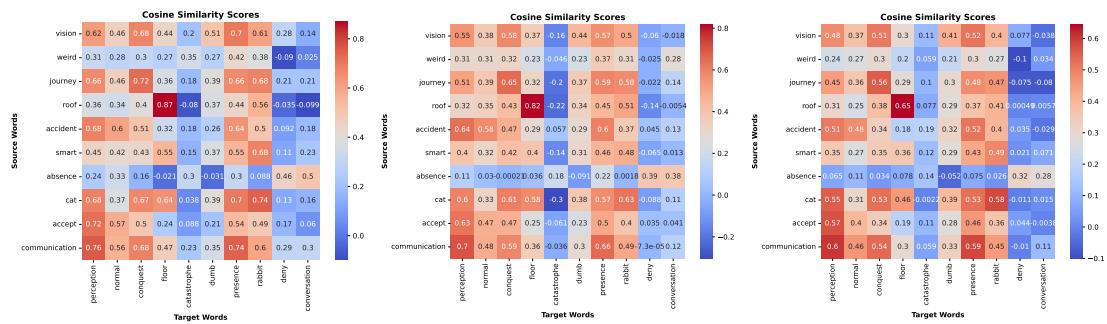


Figure 223: xho Mono Emb plots of 50 (left), 100(middle), and 200 (right), nso test

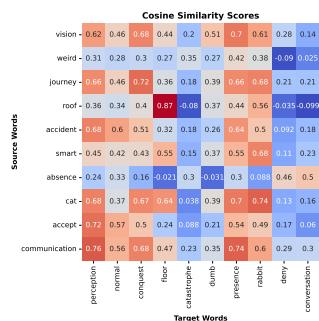


Figure 224: hr

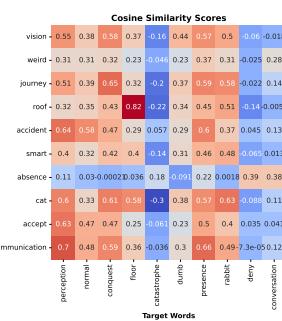


Figure 225: hr

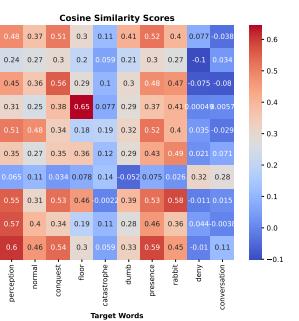


Figure 226: hr

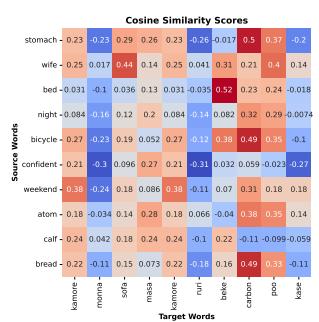


Figure 227: hr and lr-t

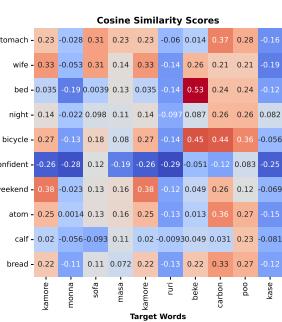


Figure 228: hr and lr-t

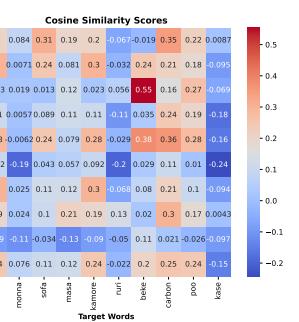


Figure 229: hr and lr-t

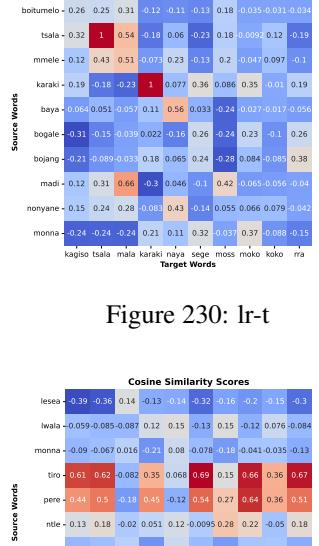
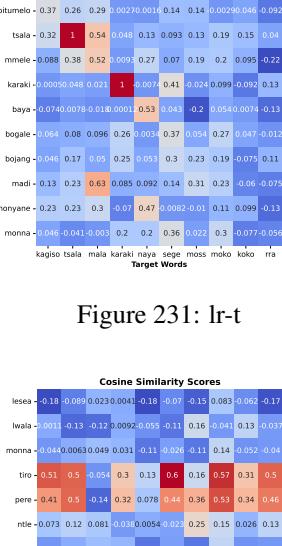


Figure 230: lr-t



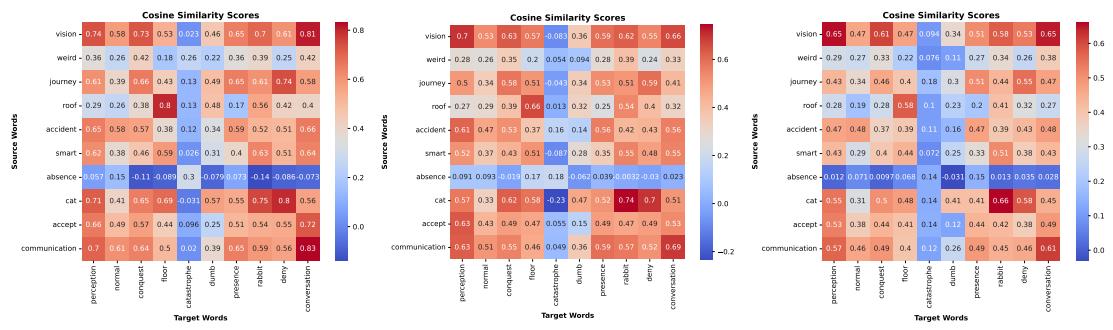


Figure 237: hr

Figure 238: hr

Figure 239: hr

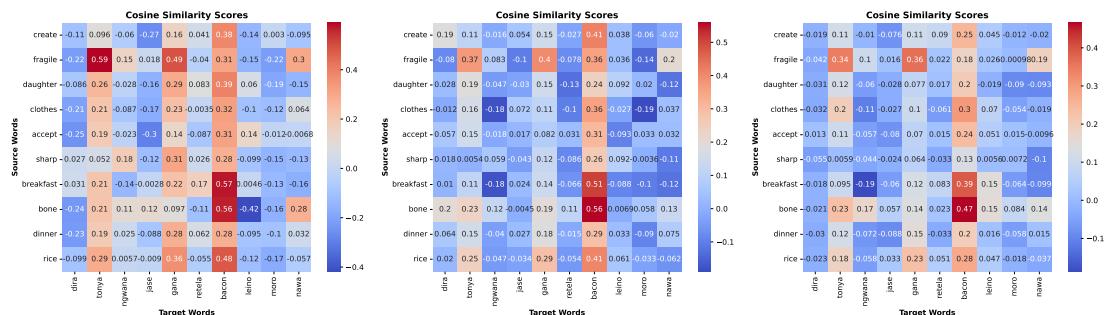


Figure 240: hr and lr-t

Figure 241: hr and lr-t

Figure 242: hr and lr-t

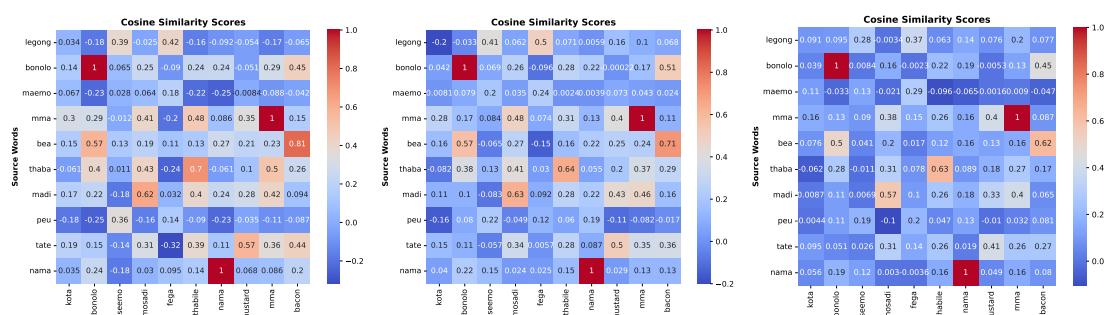


Figure 243: lr-t

Figure 244: lr-t

Figure 245: lr-t

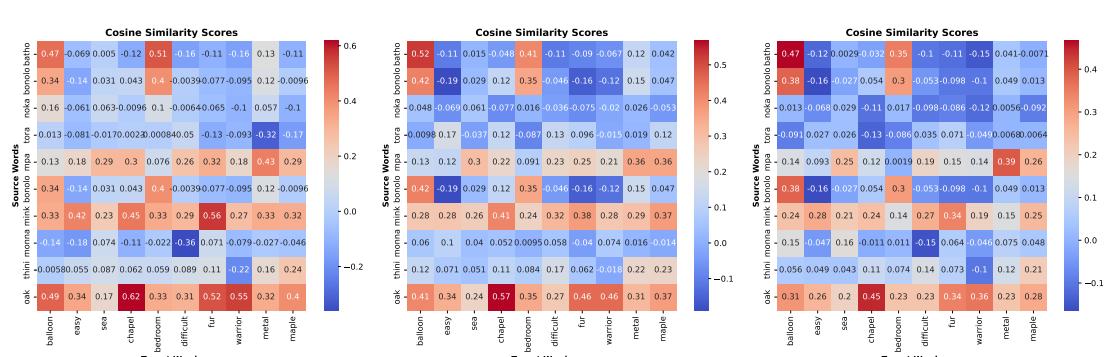


Figure 246: lr-t and hr

Figure 247: lr-t and hr

Figure 248: lr-t and hr

Figure 249: zul Mono Emb plots of 50 (left), 100(middle), and 200 (right), nso test

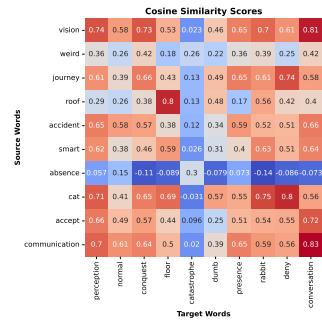


Figure 250: hr

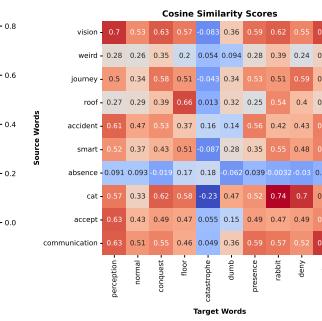


Figure 251: hr

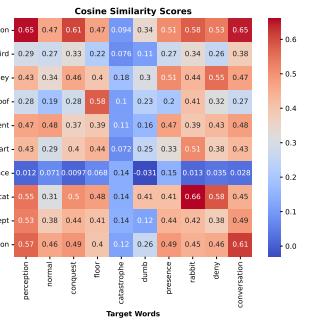


Figure 252: hr

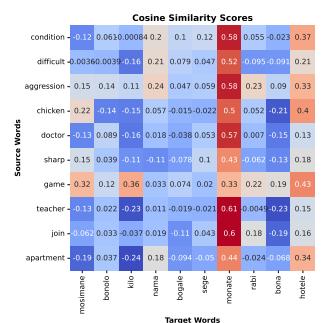


Figure 253: hr and lr-t

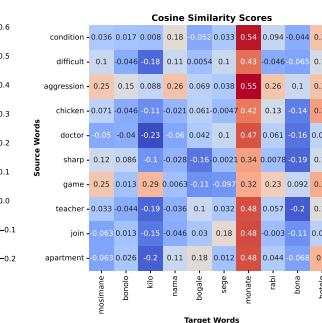


Figure 254: hr and lr-t

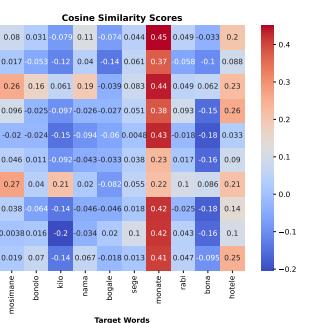


Figure 255: hr and lr-t

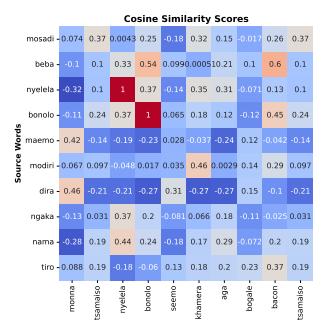


Figure 256: lr-t

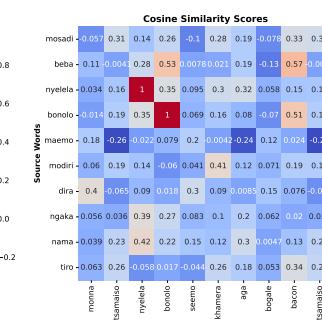


Figure 257: lr-t

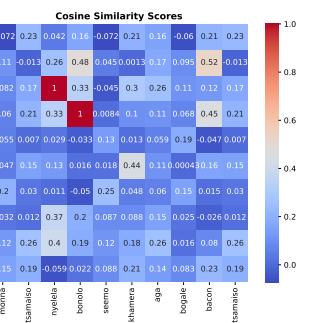


Figure 258: lr-t

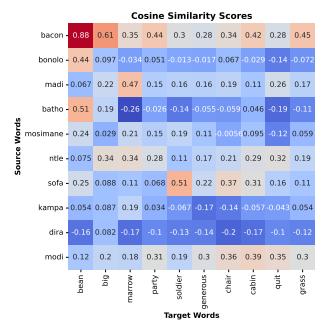


Figure 259: lr-t and hr

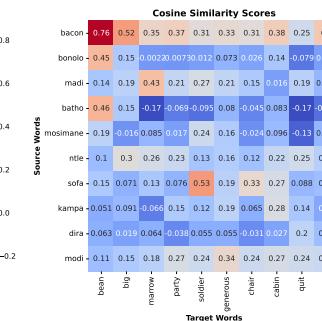


Figure 260: lr-t and hr

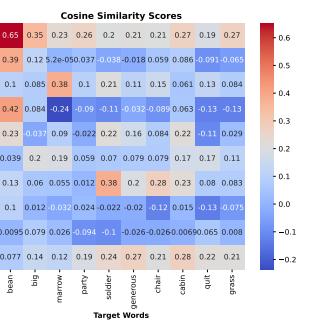
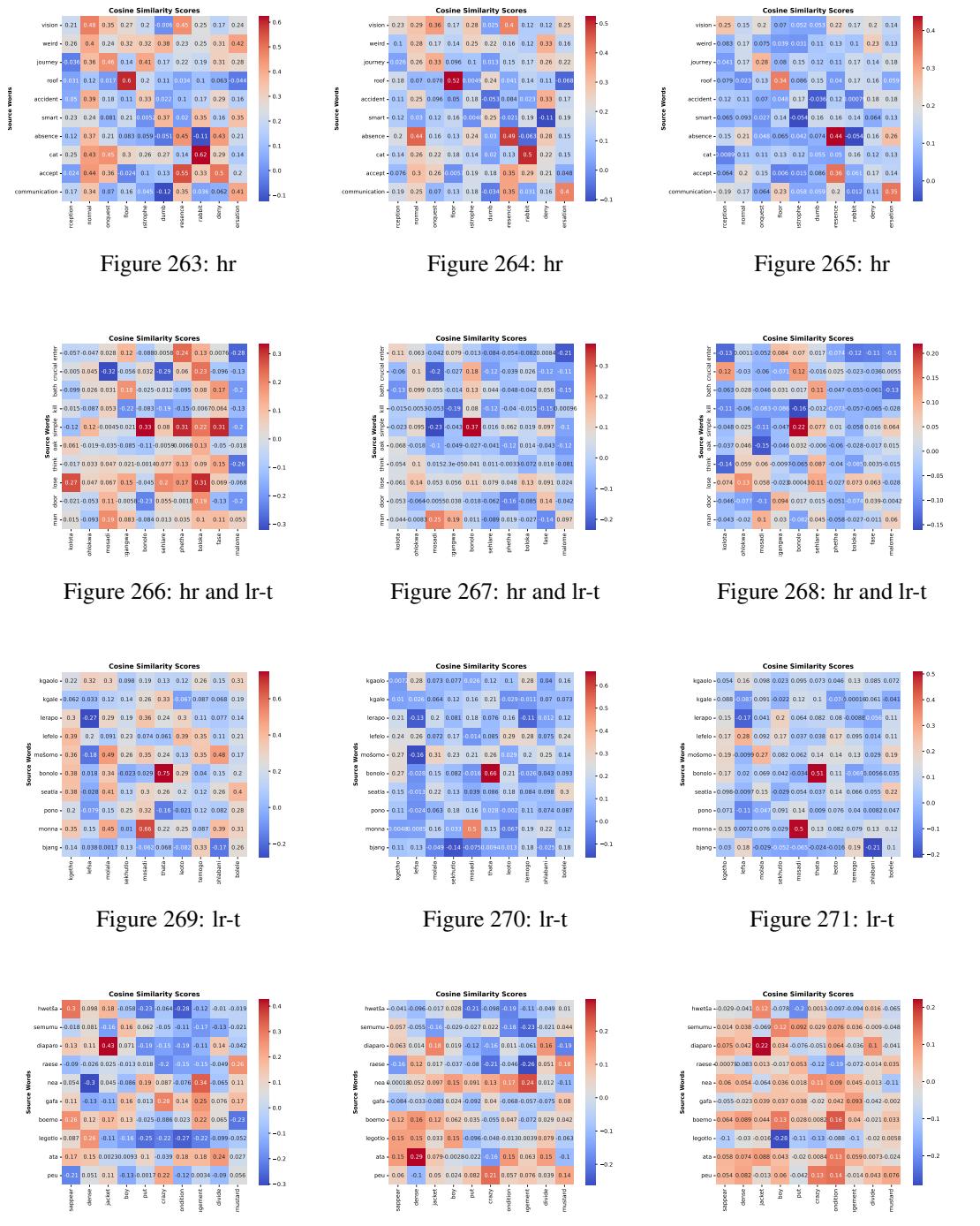


Figure 261: lr-t and hr

Figure 262: zul Mono plots of 50 (left), 100(middle), and 200 (right), tsn test



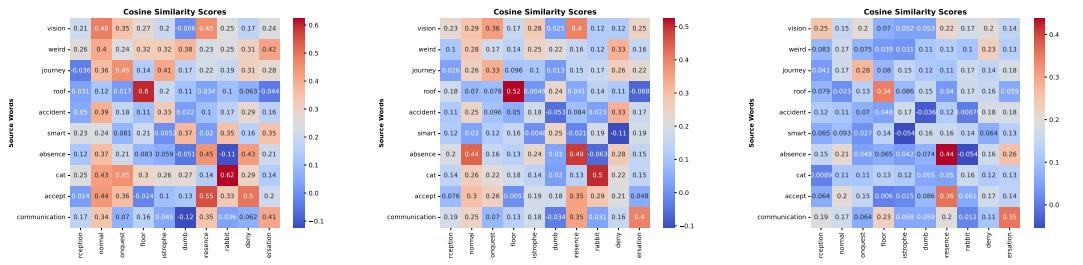


Figure 276: hr

Figure 277: hr

Figure 278: hr

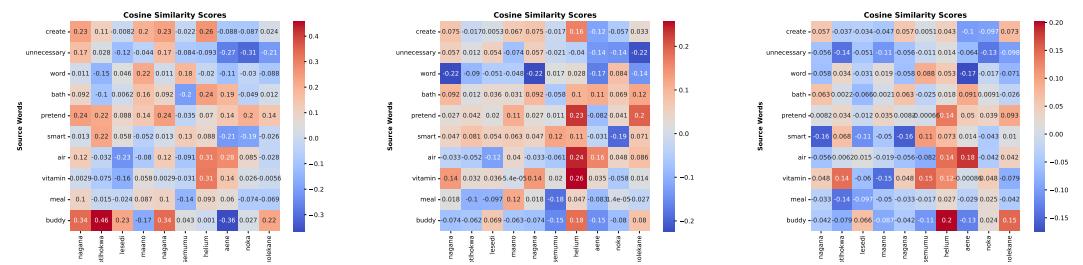


Figure 279: hr and lr-t

Figure 280: hr and lr-t

Figure 281: hr and lr-t

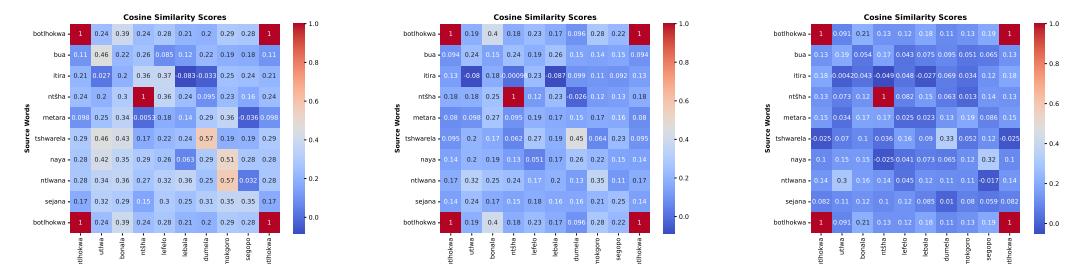


Figure 282: lr-t

Figure 283: lr-t

Figure 284: lr-t

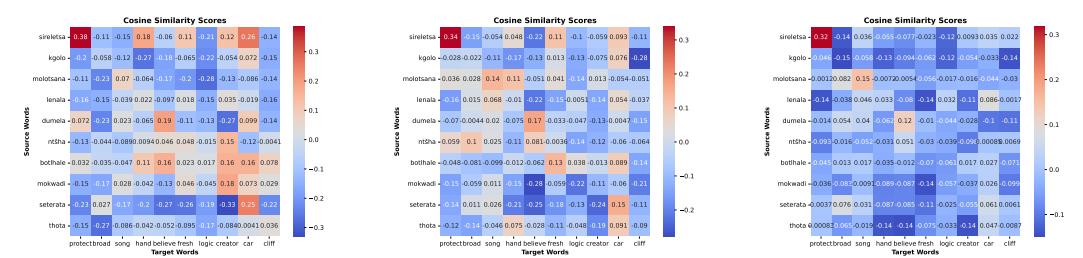


Figure 285: lr-t and hr

Figure 286: lr-t and hr

Figure 287: lr-t and hr

Figure 288: sot Emb CCA plots of 50 (left), 100(middle), and 200 (right), tsn test

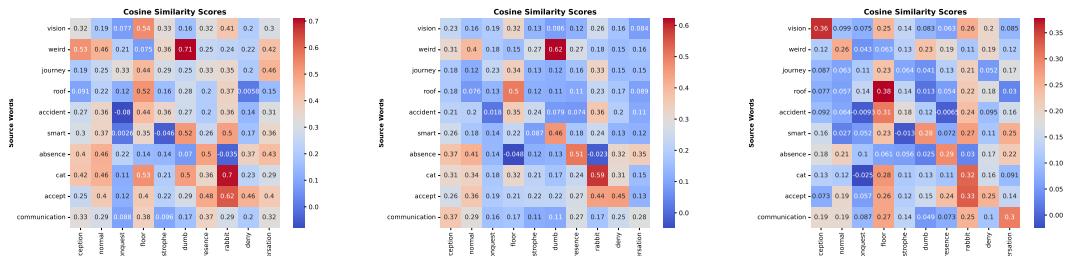


Figure 289: hr

Figure 290: hr

Figure 291: hr

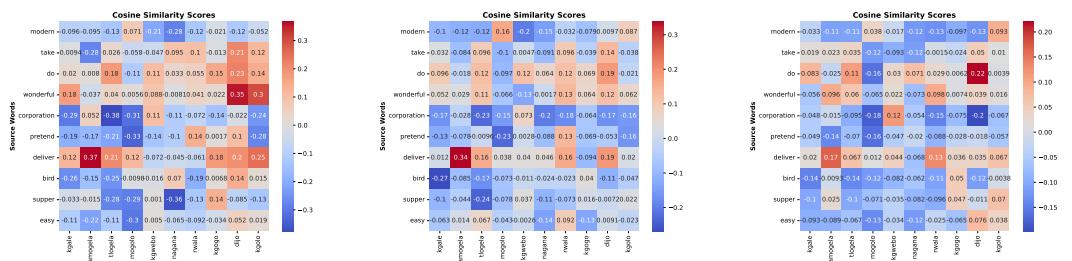


Figure 292: hr and lr-t

Figure 293: hr and lr-t

Figure 294: hr and lr-t

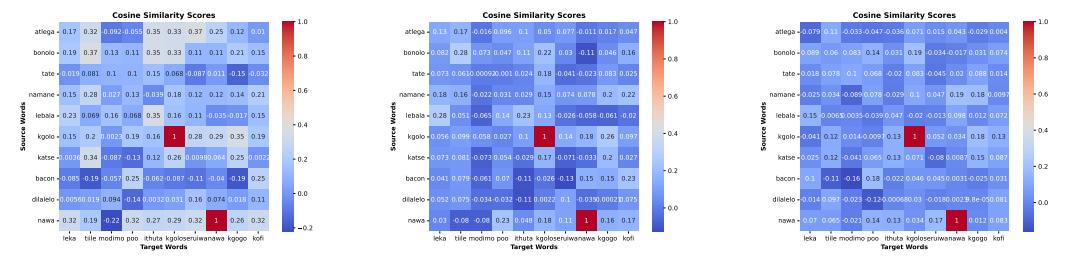


Figure 295: lr-t

Figure 296: lr-t

Figure 297: lr-t

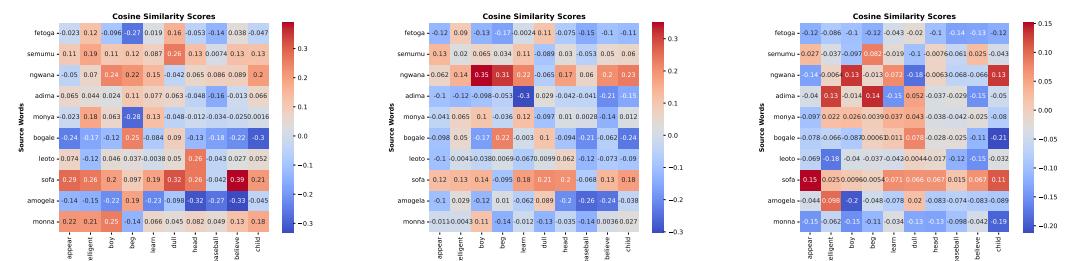


Figure 298: lr-t and hr

Figure 299: lr-t and hr

Figure 300: lr-t and hr

Figure 301: tsn CCA plots of 50 (left), 100(middle), and 200 (right), nso test

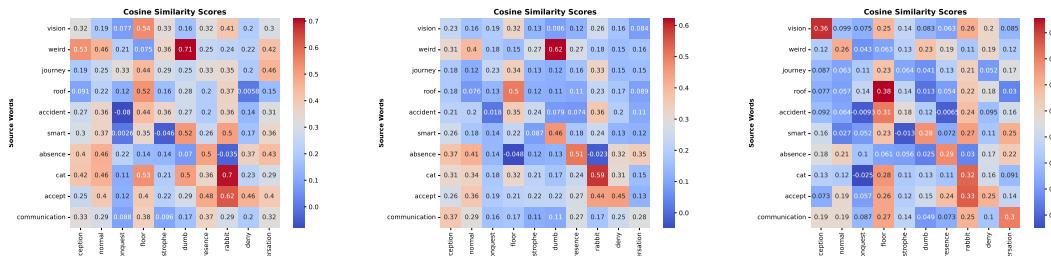


Figure 302: hr

Figure 303: hr

Figure 304: hr

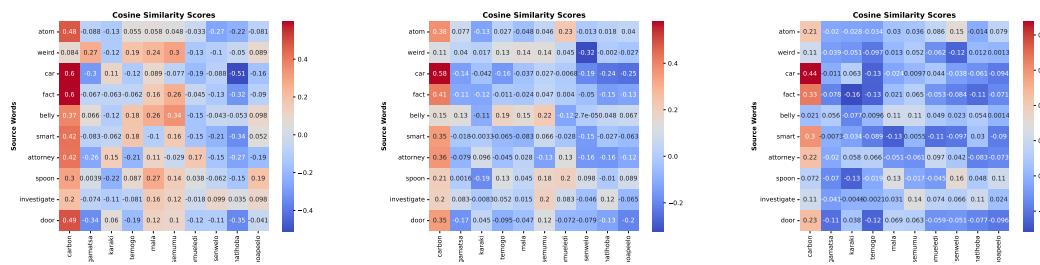


Figure 305: hr and lr-t

Figure 306: hr and lr-t

Figure 307: hr and lr-t

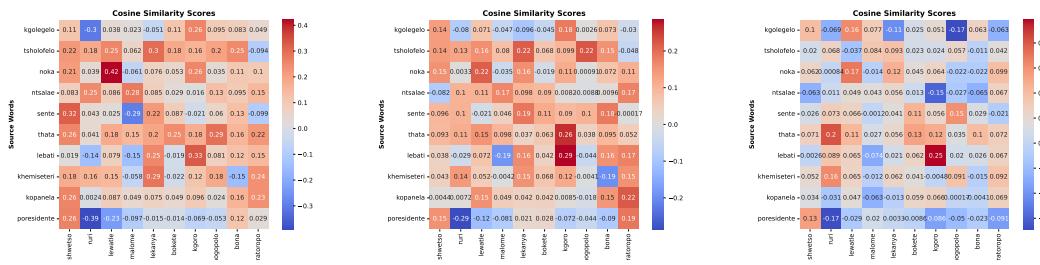


Figure 308: lr-t

Figure 309: lr-t

Figure 310: lr-t

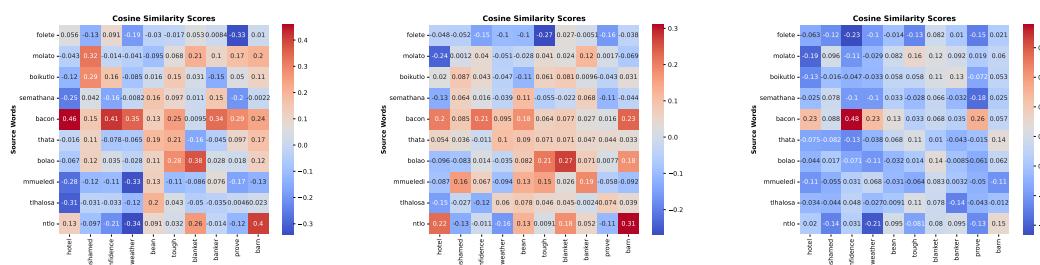


Figure 311: lr-t and hr

Figure 312: lr-t and hr

Figure 313: lr-t and hr

Figure 314: tsn CCA plots of 50 (left), 100(middle), and 200 (right), tsn test

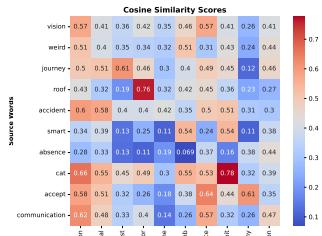


Figure 315: hr

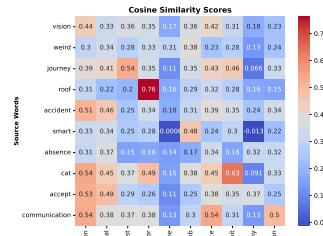


Figure 316: hr

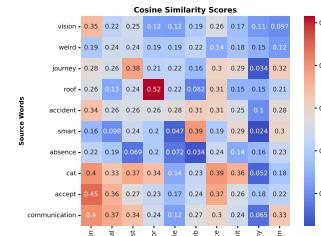


Figure 317: hr

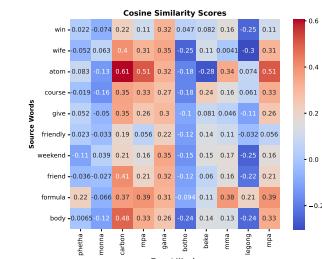


Figure 318: hr and lr-t

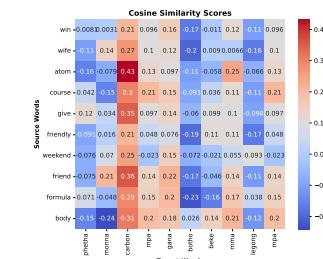


Figure 319: hr and lr-t

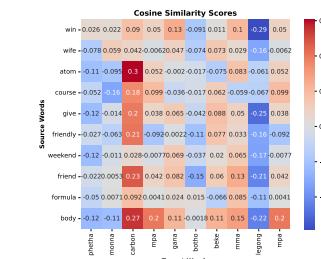


Figure 320: hr and lr-t

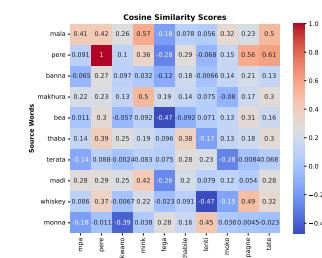


Figure 321: lr-t

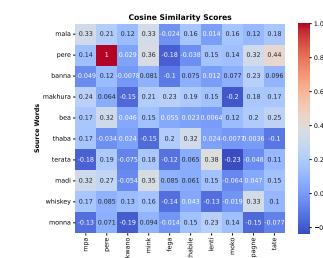


Figure 322: lr-t

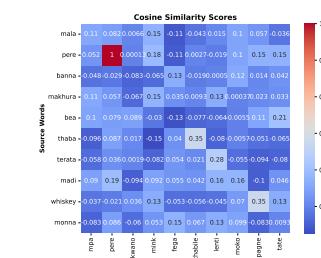


Figure 323: lr-t

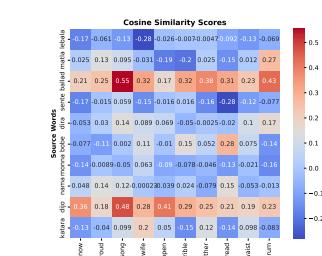


Figure 324: lr-t and hr

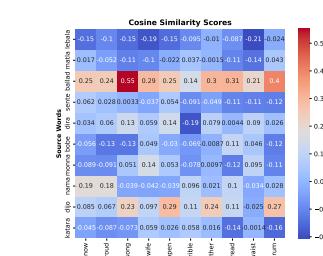


Figure 325: lr-t and hr

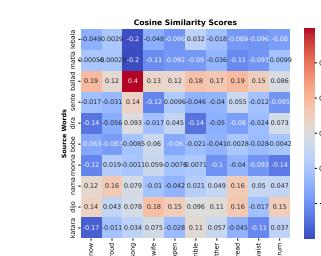


Figure 326: lr-t and hr

Figure 327: xho CCA plots of 50 (left), 100(middle), and 200 (right), nso test

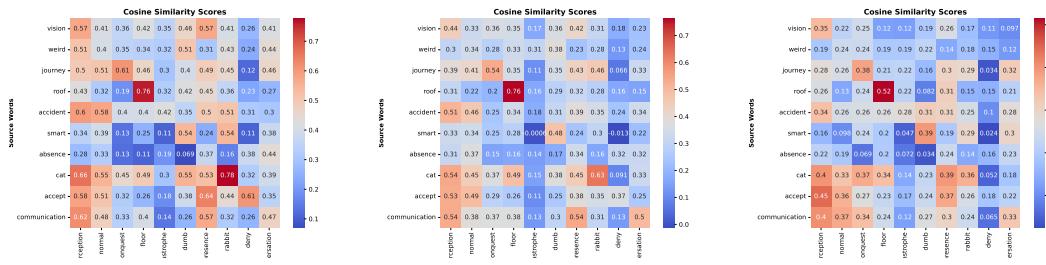


Figure 328: hr

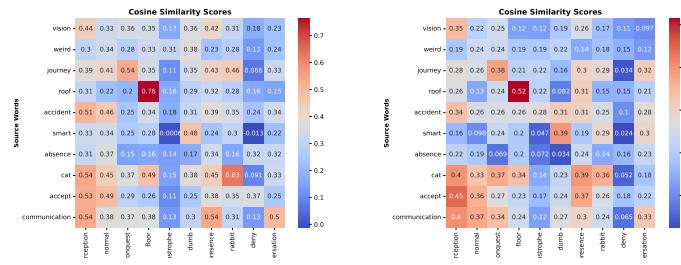


Figure 329: hr

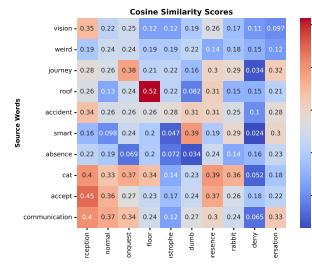


Figure 330: hr

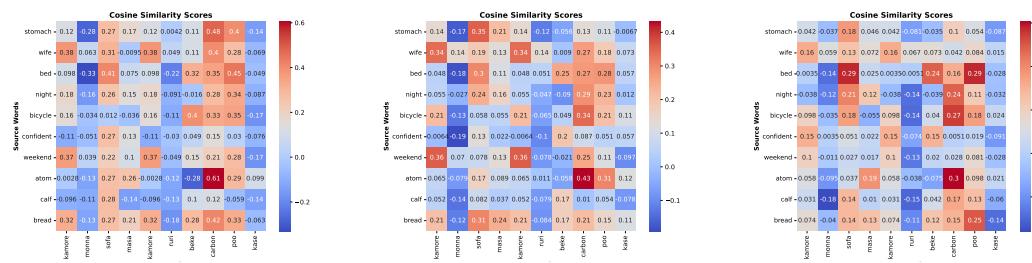


Figure 331: hr and lr-t

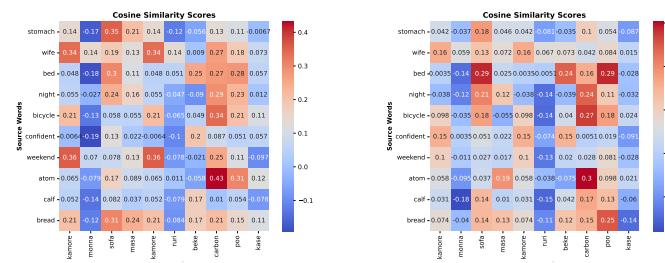


Figure 332: hr and lr-t

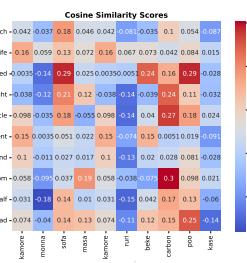


Figure 333: hr and lr-t

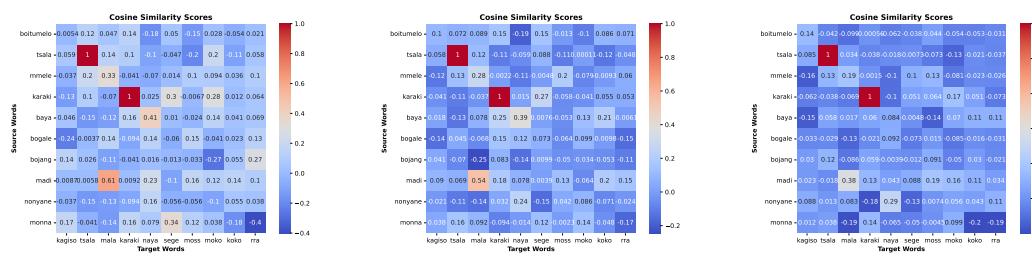


Figure 334: lr-t

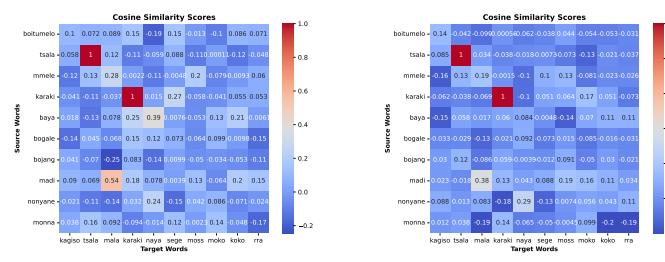


Figure 335: lr-t

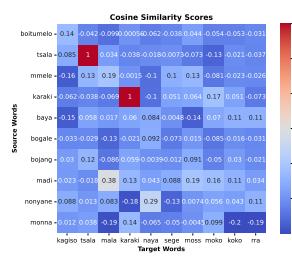


Figure 336: lr-t

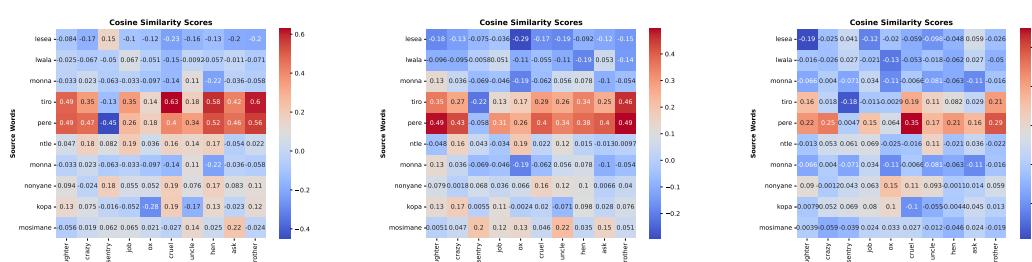


Figure 337: lr-t and hr

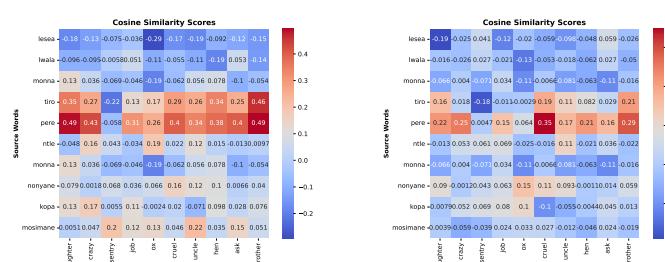


Figure 338: lr-t and hr

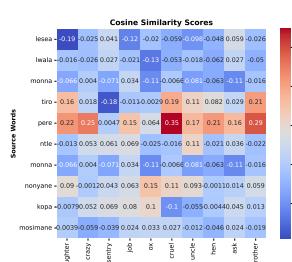


Figure 339: lr-t and hr

Figure 340: xho CCA plots of 50 (left), 100(middle), and 200 (right), tsn test

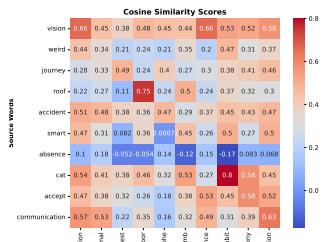


Figure 341: hr

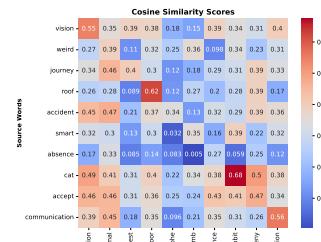


Figure 342: hr

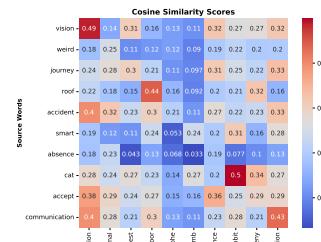


Figure 343: hr

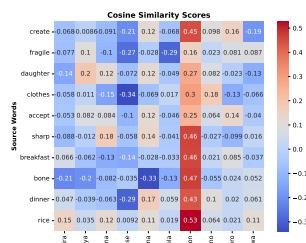


Figure 344: hr and lr-t

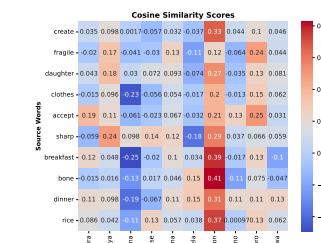


Figure 345: hr and lr-t

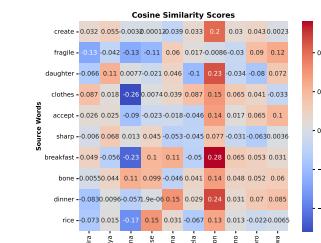


Figure 346: hr and lr-t

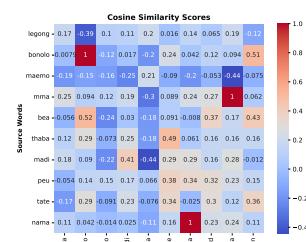


Figure 347: lr-t

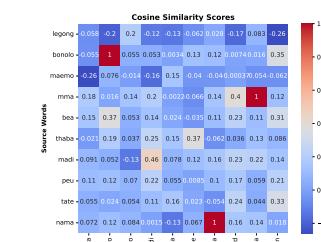


Figure 348: lr-t

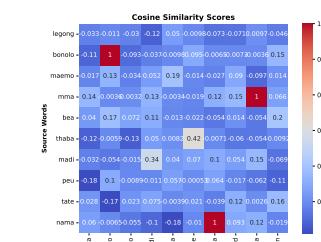


Figure 349: lr-t

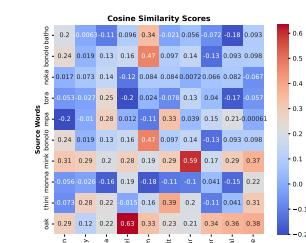


Figure 350: lr-t and hr

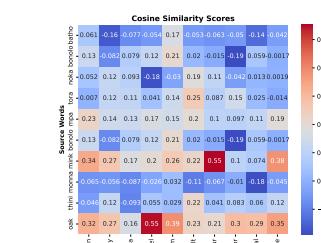


Figure 351: lr-t and hr

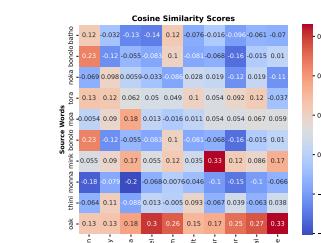


Figure 352: lr-t and hr

Figure 353: zul CCA plots of 50 (left), 100(middle), and 200 (right), nso test

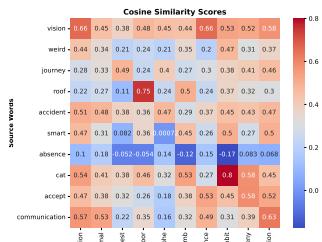


Figure 354: hr

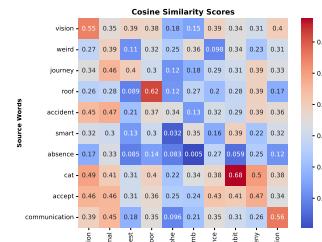


Figure 355: hr

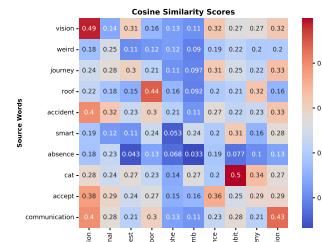


Figure 356: hr

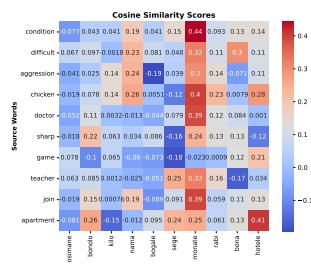


Figure 357: hr and lr-t

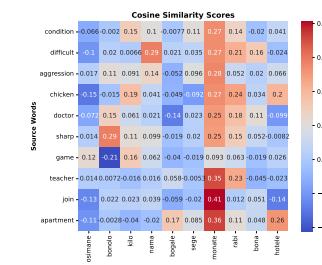


Figure 358: hr and lr-t

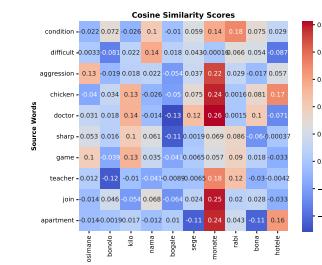


Figure 359: hr and lr-t

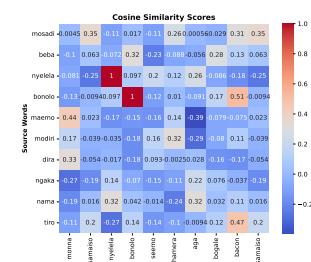


Figure 360: lr-t

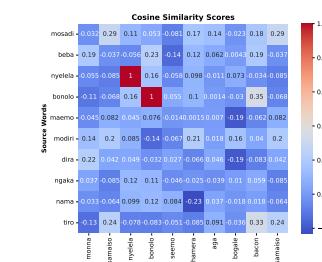


Figure 361: lr-t

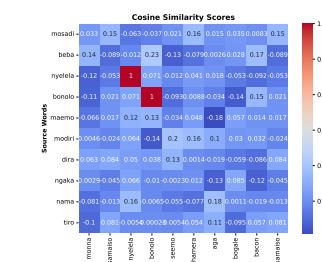


Figure 362: lr-t

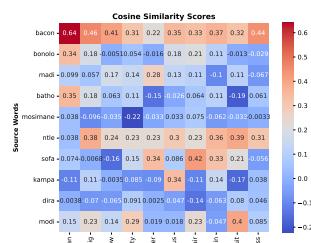


Figure 363: lr-t and hr

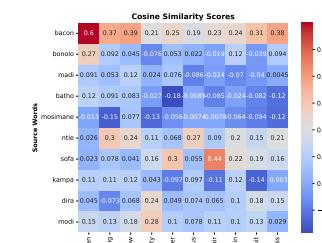


Figure 364: lr-t and hr

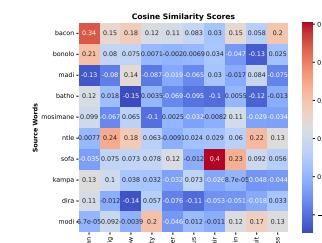


Figure 365: lr-t and hr

Figure 366: zul CCA plots of 50 (left), 100(middle), and 200 (right), tsn test

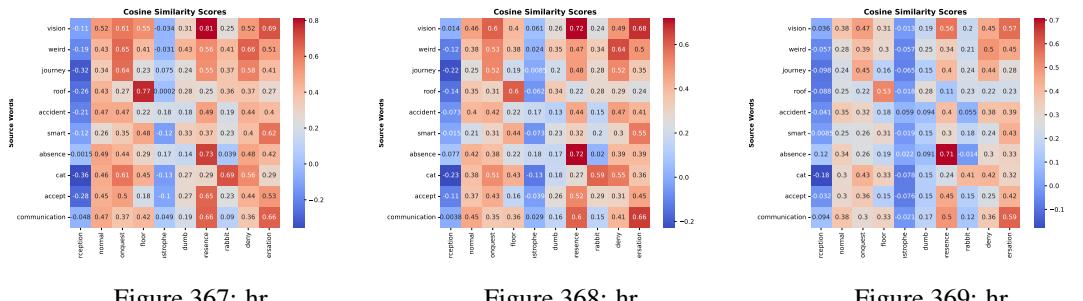


Figure 367: hr

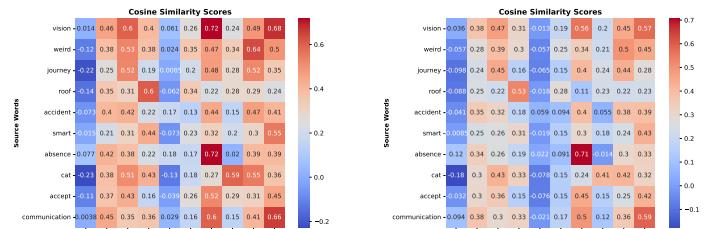


Figure 368: hr

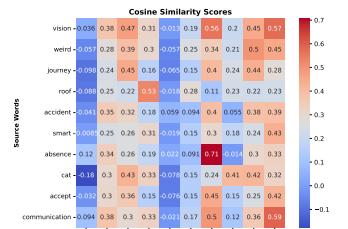


Figure 369: hr

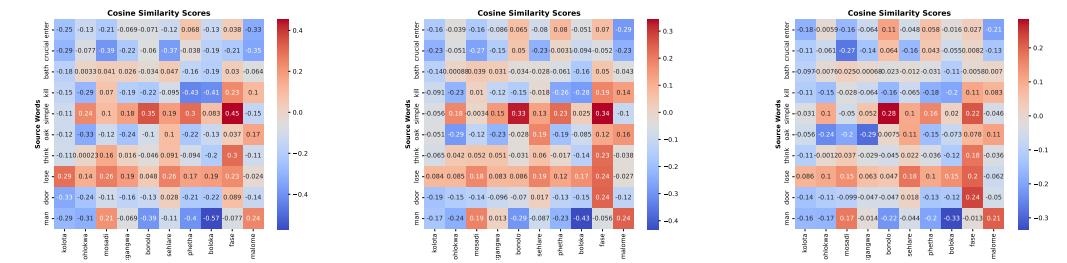


Figure 370: hr and lr-t

Figure 371: hr and lr-t

Figure 372: hr and lr-t

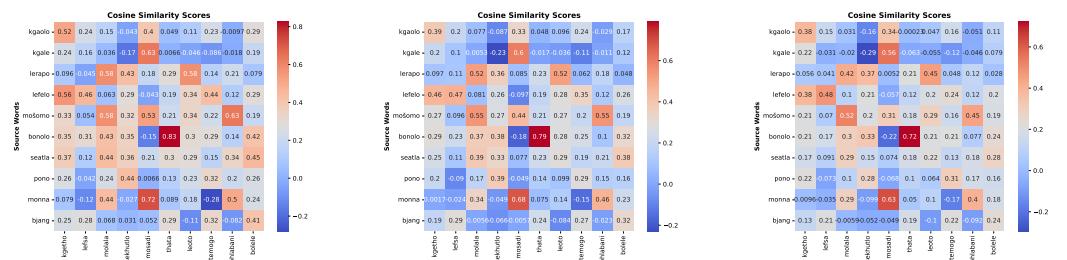


Figure 373: lr-t

Figure 374: lr-t

Figure 375: lr-t

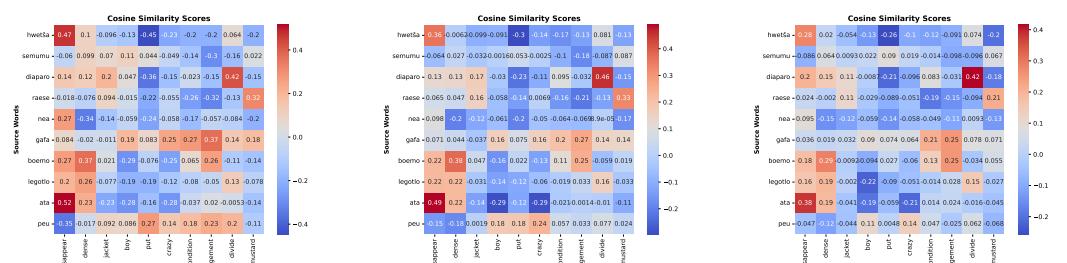


Figure 376: lr-t and hr

Figure 377: lr-t and hr

Figure 378: lr-t and hr

Figure 379: sot Emb Muse plots of 50 (left), 100(middle), and 200 (right), nso test

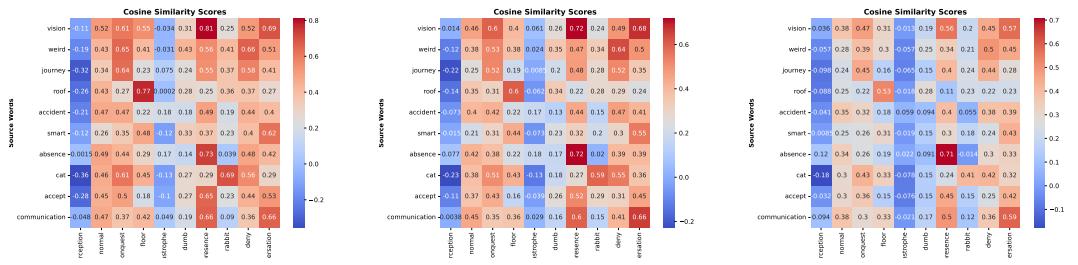


Figure 380: hr

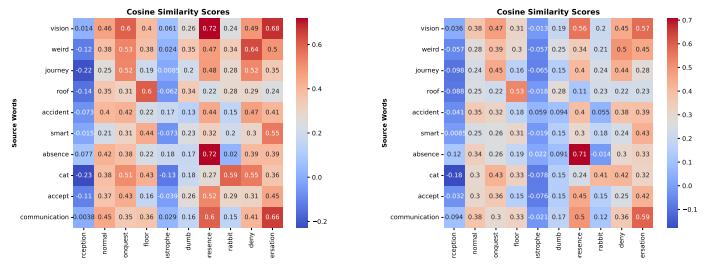


Figure 381: hr

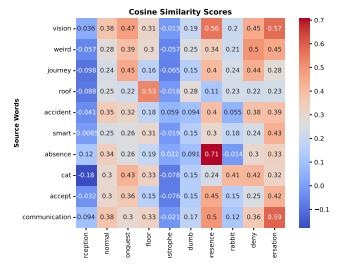


Figure 382: hr

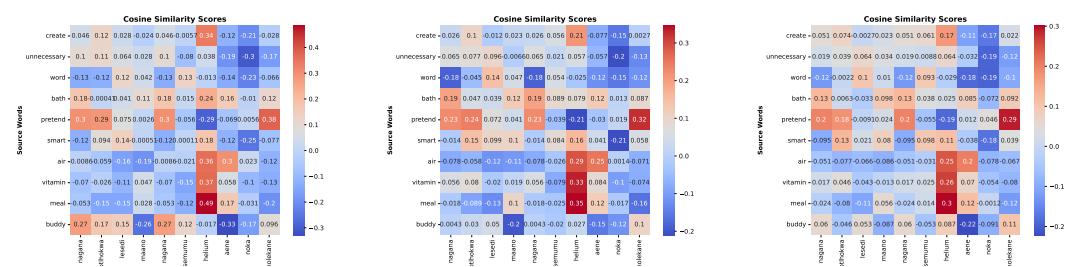


Figure 383: hr and lr-t

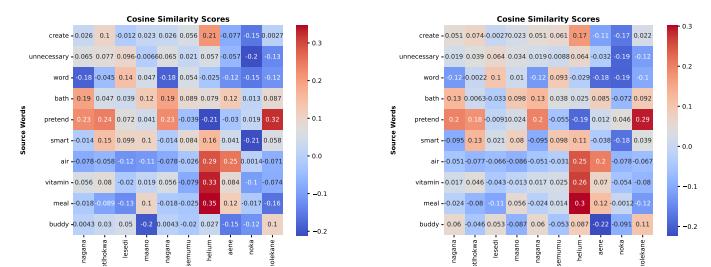


Figure 384: hr and lr-t

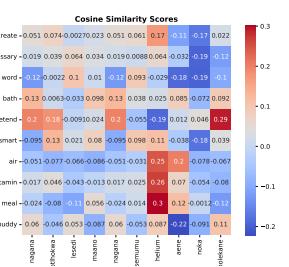


Figure 385: hr and lr-t

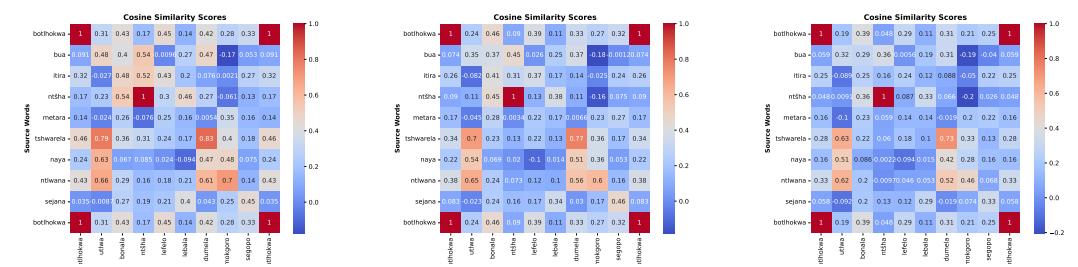


Figure 386: lr-t

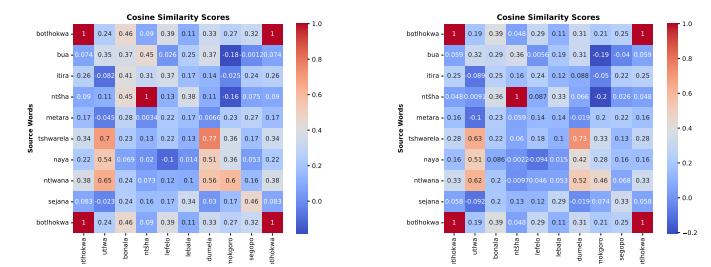


Figure 387: lr-t

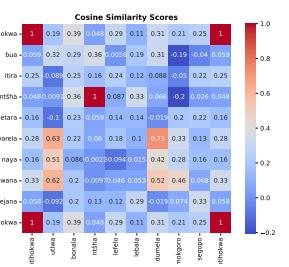


Figure 388: lr-t

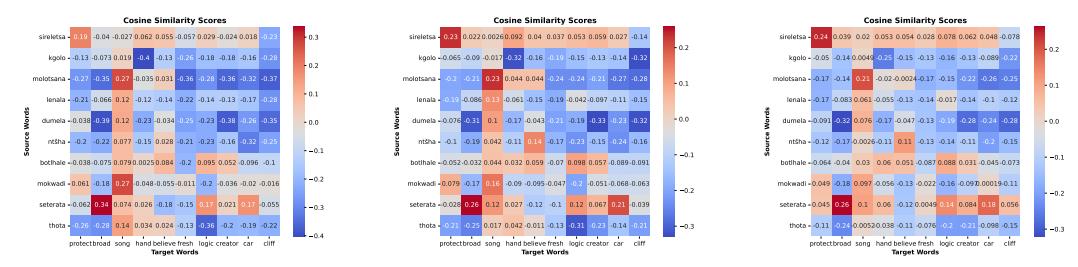


Figure 389: lr-t and hr

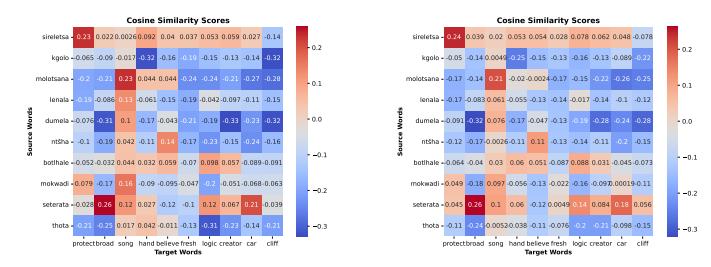


Figure 390: lr-t and hr

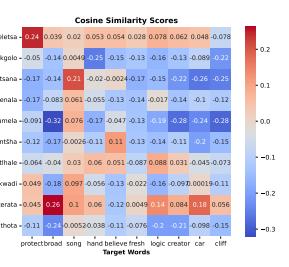


Figure 391: lr-t and hr

Figure 392: sot Emb Muse plots of 50 (left), 100(middle), and 200 (right), tsn test

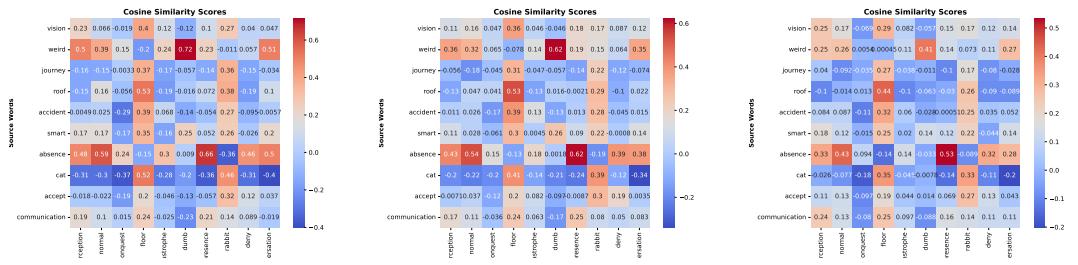


Figure 393: hr

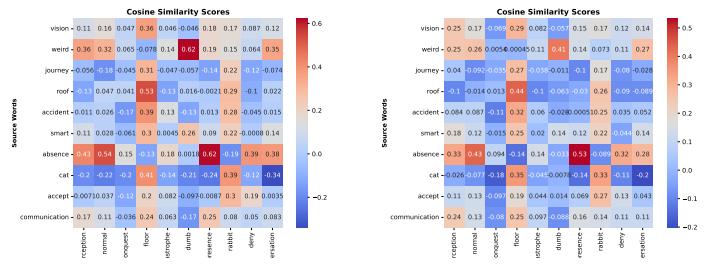


Figure 394: hr

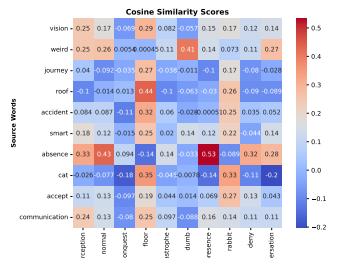


Figure 395: hr

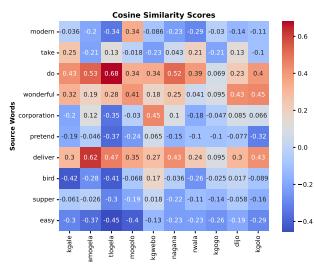


Figure 396: hr and lr-t

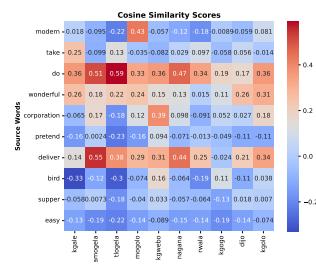


Figure 397: hr and lr-t

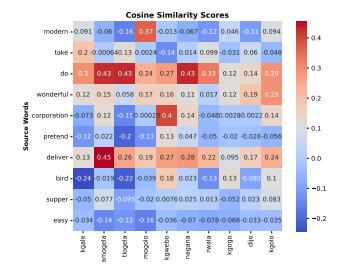


Figure 398: hr and lr-t

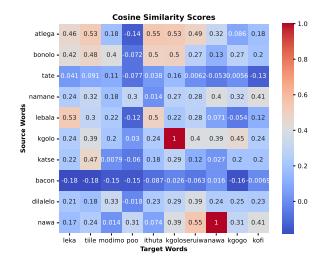


Figure 399: lr-t

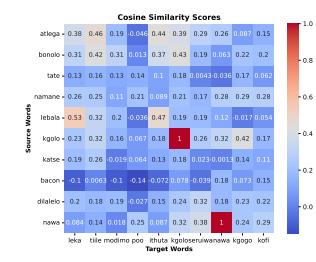


Figure 400: lr-t

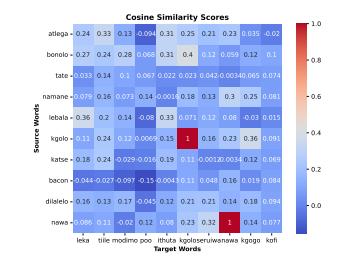


Figure 401: lr-t

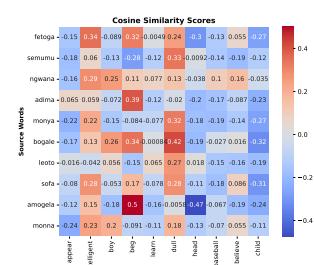


Figure 402: lr-t and hr

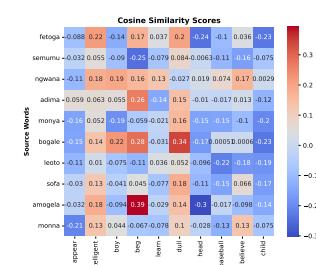


Figure 403: lr-t and hr

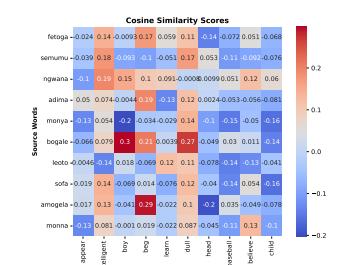


Figure 404: lr-t and hr

Figure 405: tsn Muse plots of 50 (left), 100(middle), and 200 (right), nso test

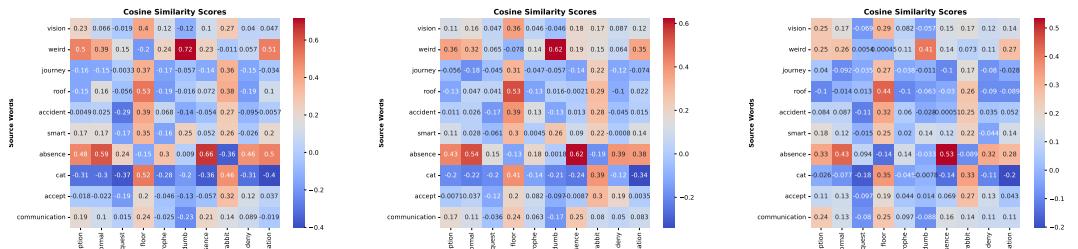


Figure 406: hr

Figure 407: hr

Figure 408: hr

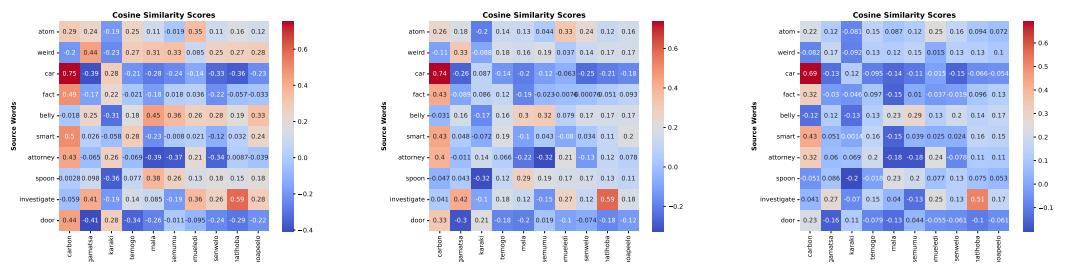


Figure 409: hr and lr-t

Figure 410: hr and lr-t

Figure 411: hr and lr-t

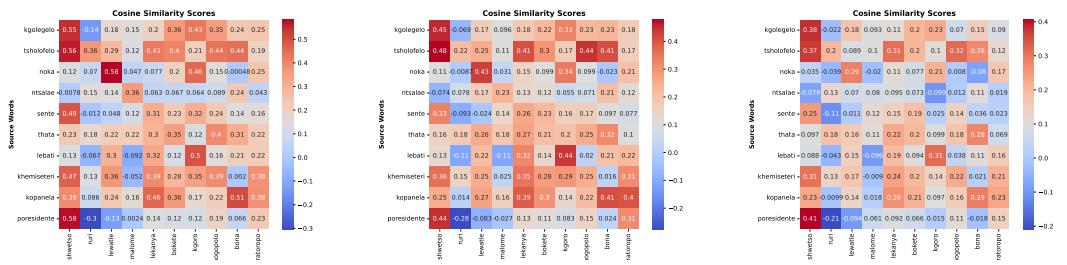


Figure 412: lr-t

Figure 413: lr-t

Figure 414: lr-t

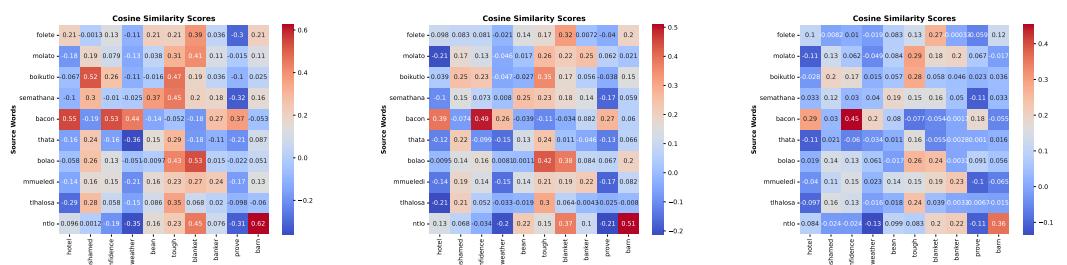


Figure 415: lr-t and hr

Figure 416: lr-t and hr

Figure 417: lr-t and hr

Figure 418: tsn Muse plots of 50 (left), 100(middle), and 200 (right), tsn test

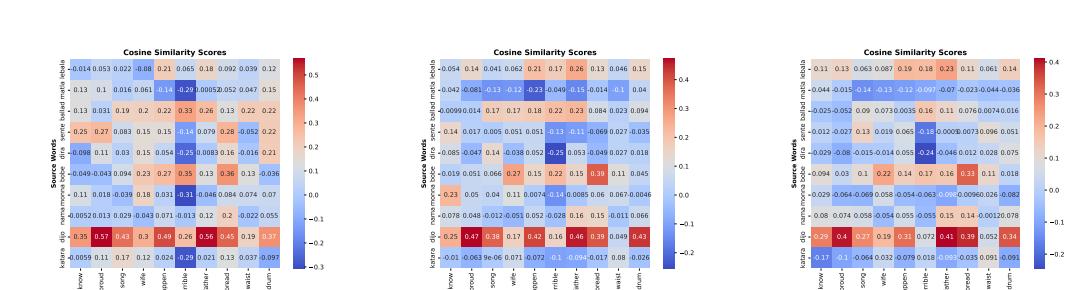
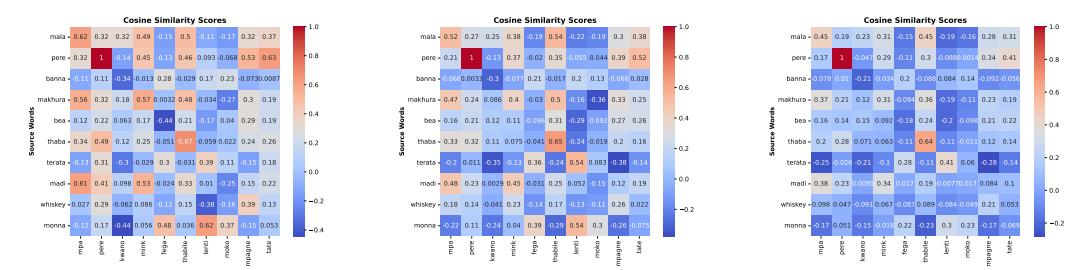
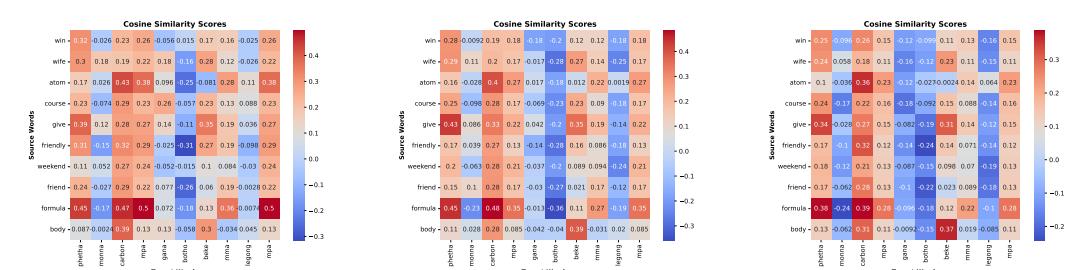
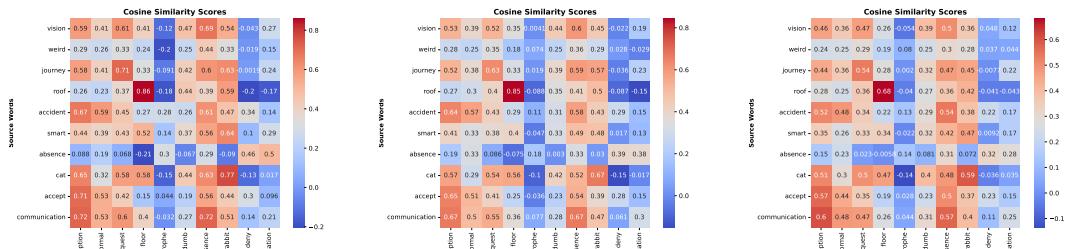


Figure 431: xho Muse plots of 50 (left), 100(middle), and 200 (right), nso test

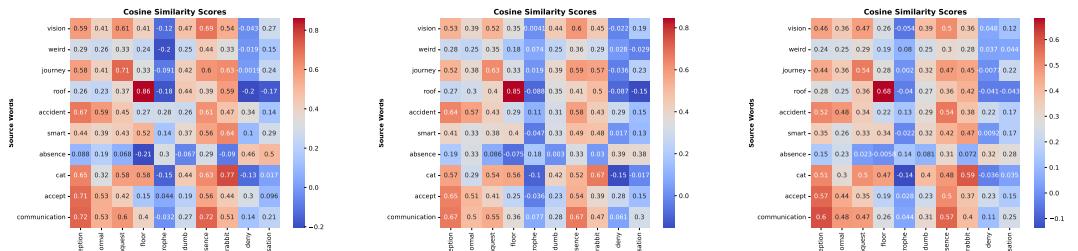


Figure 432: hr

Figure 433: hr

Figure 434: hr

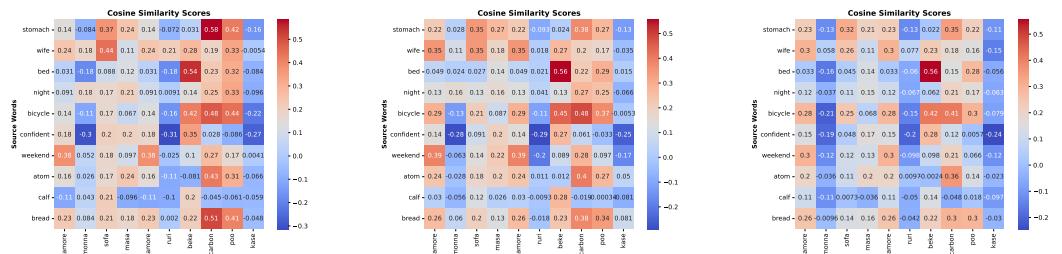


Figure 435: hr and lr-t

Figure 436: hr and lr-t

Figure 437: hr and lr-t

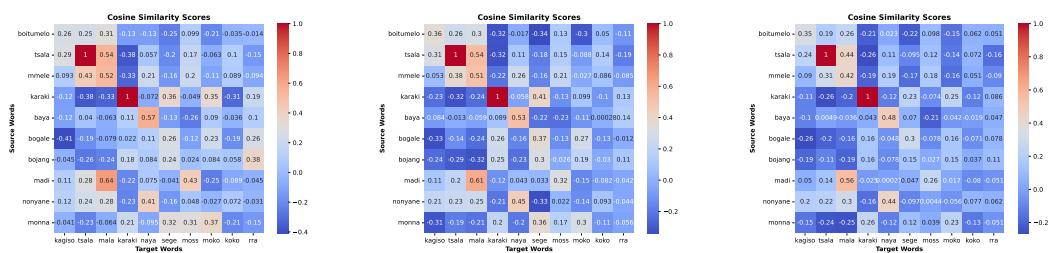


Figure 438: lr-t

Figure 439: lr-t

Figure 440: lr-t

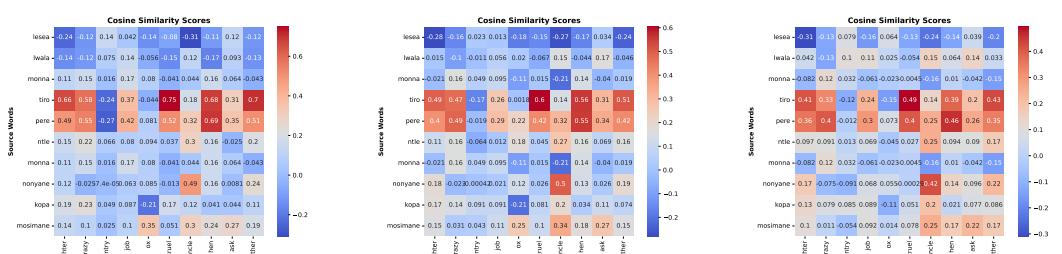


Figure 441: lr-t and hr

Figure 442: lr-t and hr

Figure 443: lr-t and hr

Figure 444: xho Muse plots of 50 (left), 100(middle), and 200 (right), tsn test

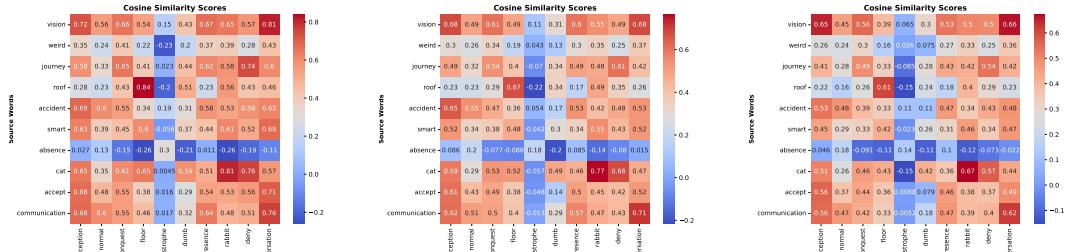


Figure 445: hr

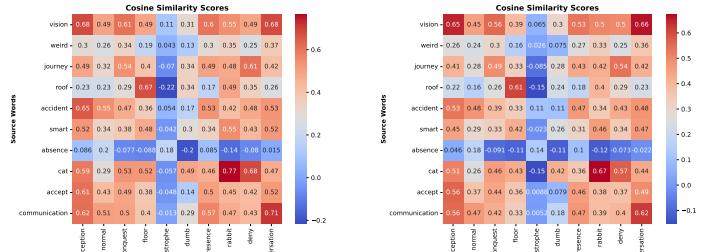


Figure 446: hr

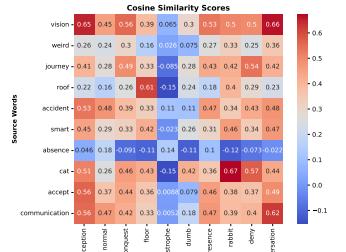


Figure 447: hr

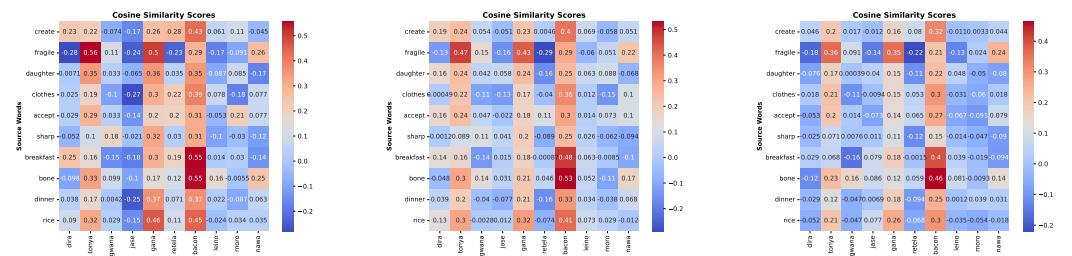


Figure 448: hr and lr-t

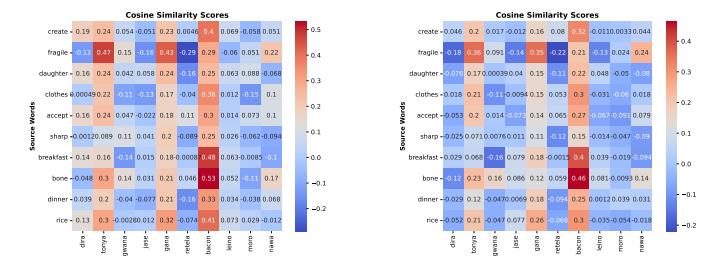


Figure 449: hr and lr-t

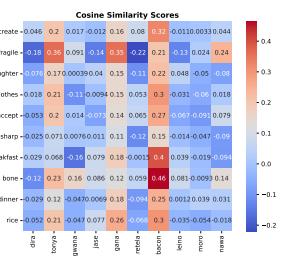


Figure 450: hr and lr-t

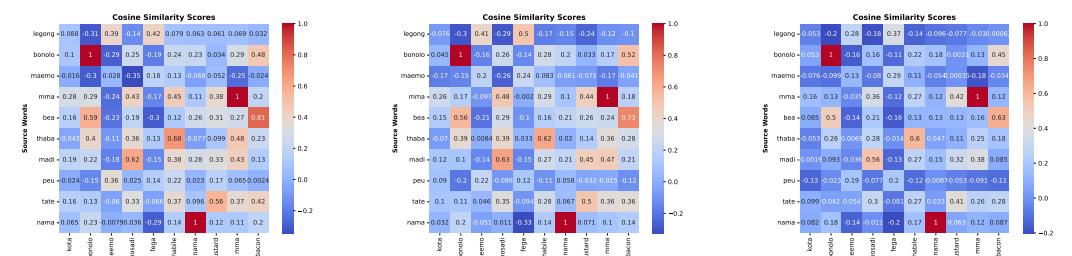


Figure 451: lr-t

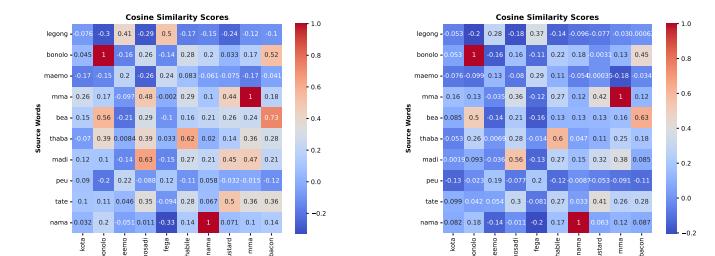


Figure 452: lr-t

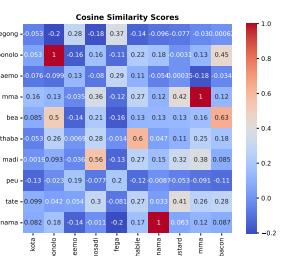


Figure 453: lr-t

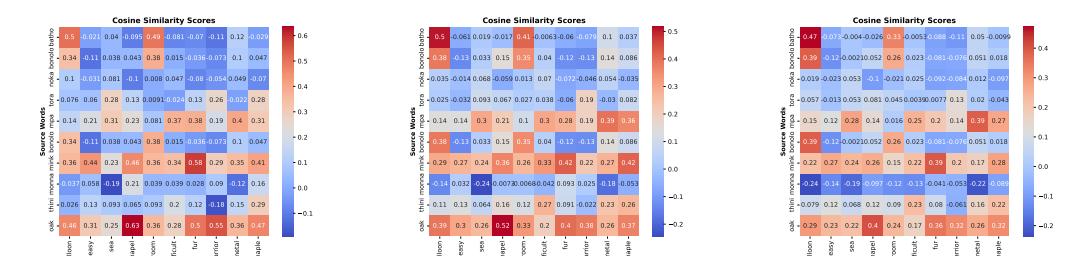


Figure 454: lr-t and hr

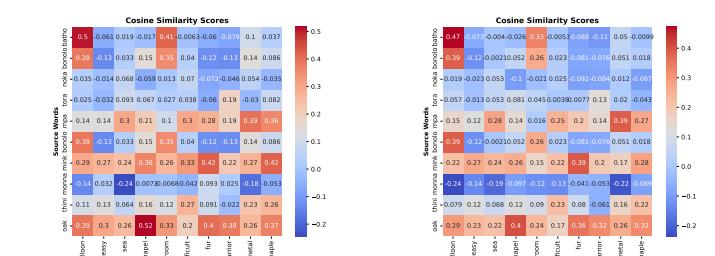


Figure 455: lr-t and hr

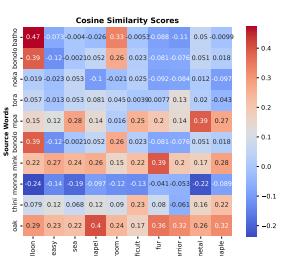


Figure 456: lr-t and hr

Figure 457: zul Muse plots of 50 (left), 100(middle), and 200 (right), nso test

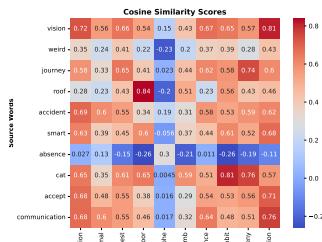


Figure 458: hr

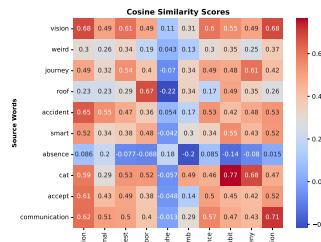


Figure 459: hr

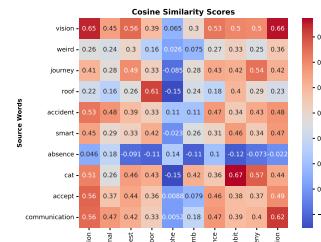


Figure 460: hr

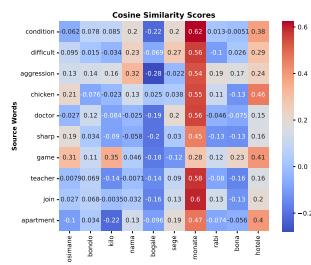


Figure 461: hr and lr-t

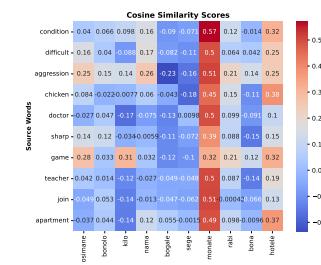


Figure 462: hr and lr-t

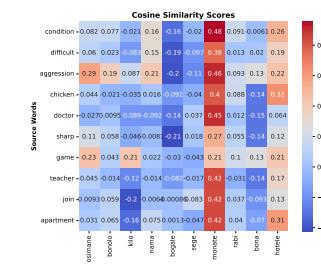


Figure 463: hr and lr-t

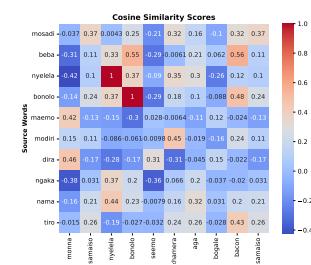


Figure 464: lr-t

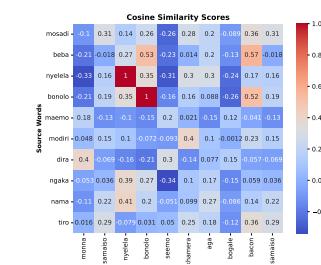


Figure 465: lr-t

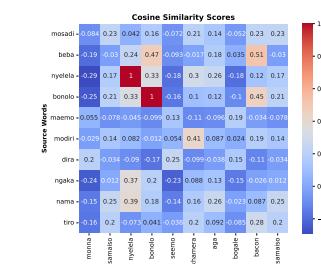


Figure 466: lr-t

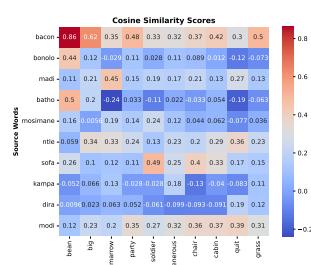


Figure 467: lr-t and hr

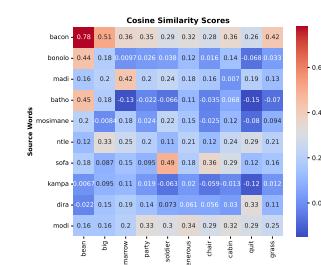


Figure 468: lr-t and hr

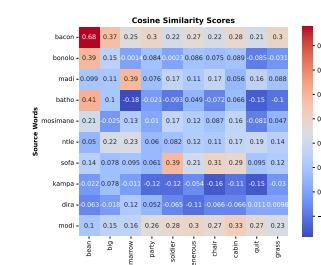


Figure 469: lr-t and hr

Figure 470: zul VecMap plots of 50 (left), 100(middle), and 200 (right), tsn test

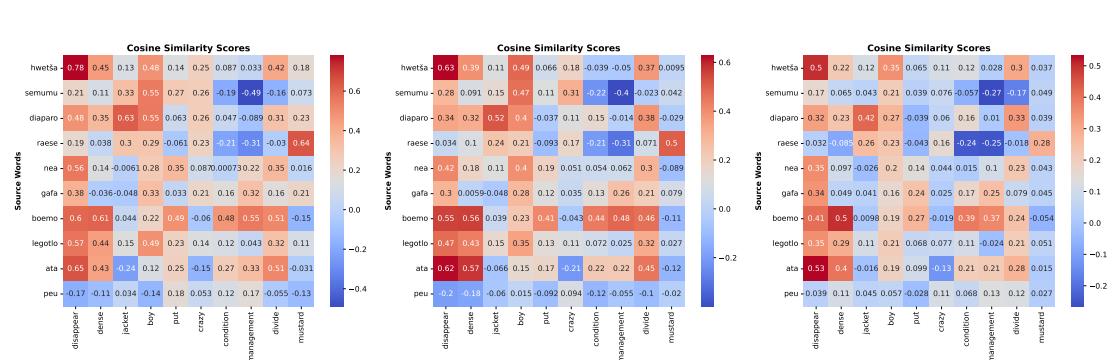
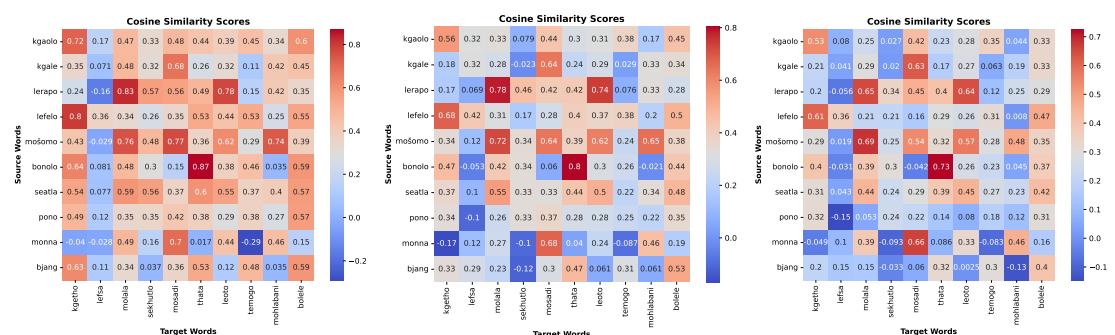
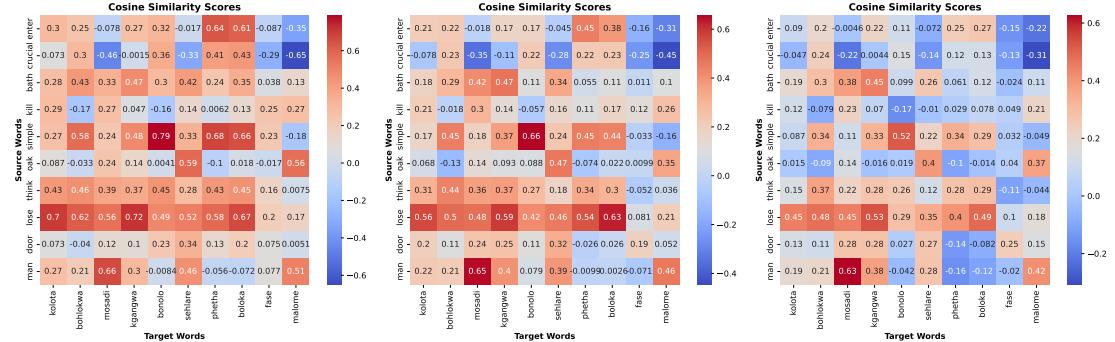
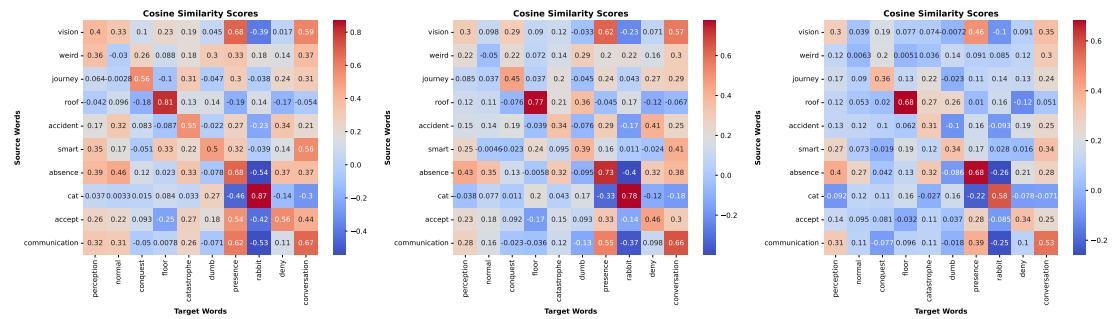


Figure 483: sot Emb VecMap plots of 50 (left), 100(middle), and 200 (right), nso test

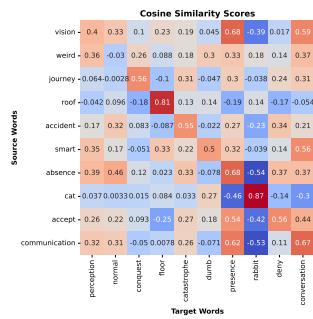


Figure 484: hr

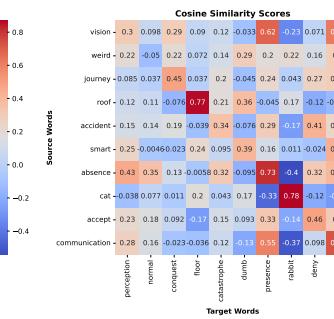


Figure 485: hr

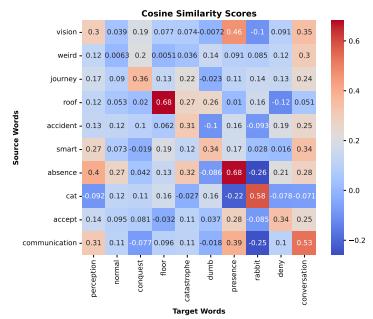


Figure 486: hr

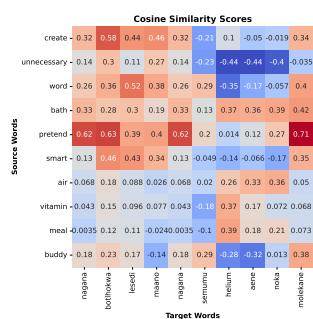


Figure 487: hr and lr-t

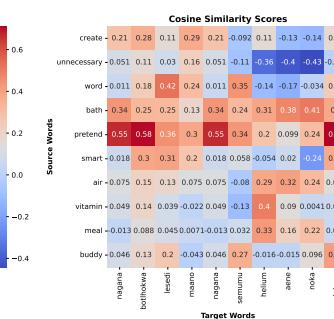


Figure 488: hr and lr-t

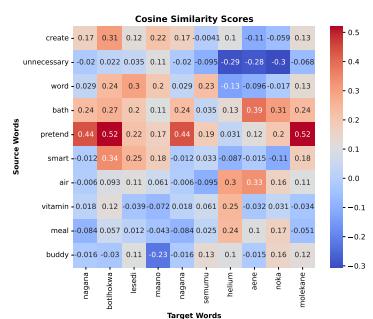


Figure 489: hr and lr-t

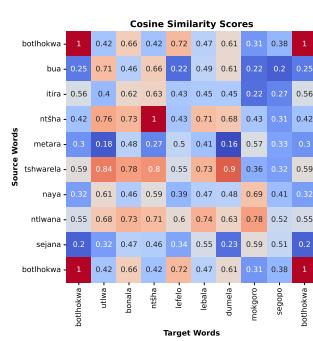


Figure 490: lr-t

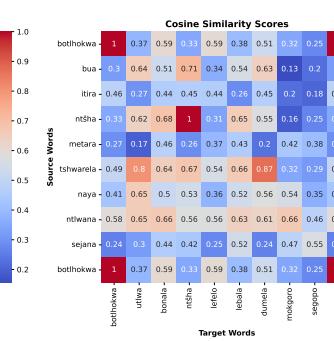


Figure 491: lr-t

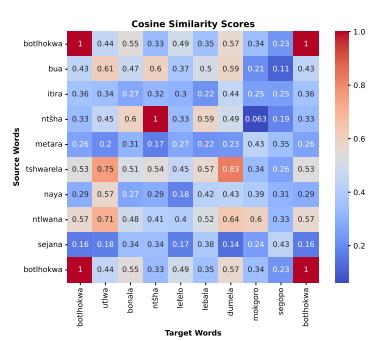


Figure 492: lr-t

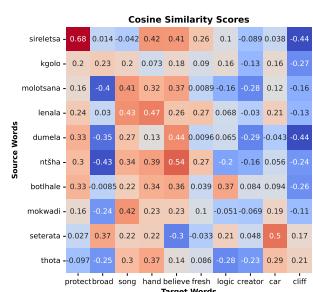


Figure 493: lr-t and hr

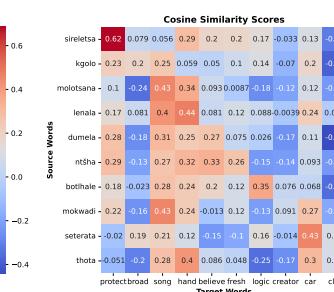


Figure 494: lr-t and hr

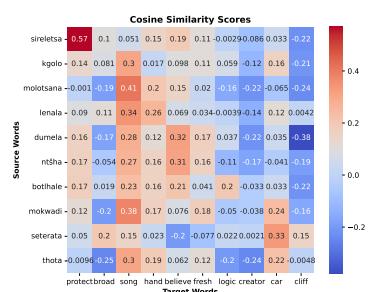


Figure 495: lr-t and hr

Figure 496: sot Emb VecMap plots of 50 (left), 100(middle), and 200 (right), tsn test

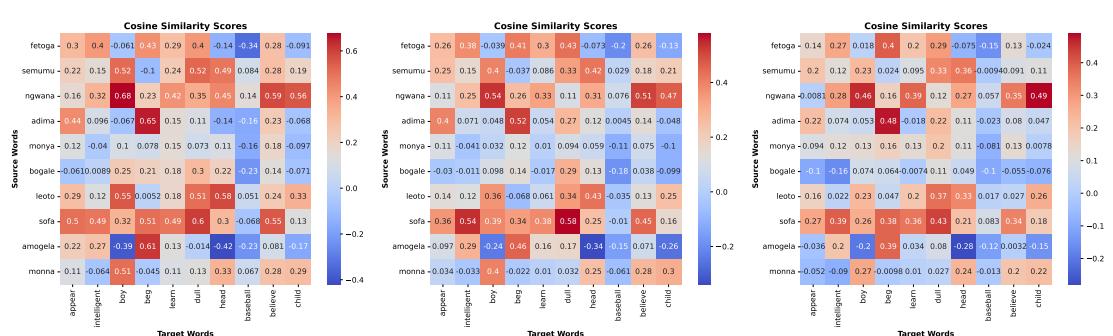
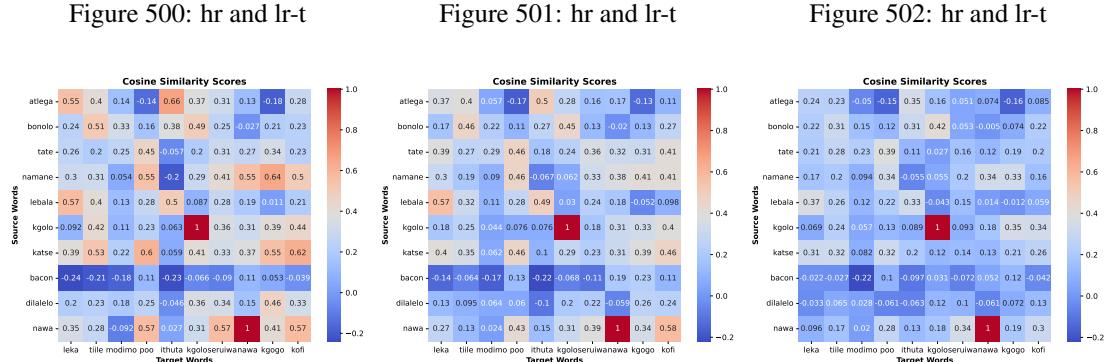
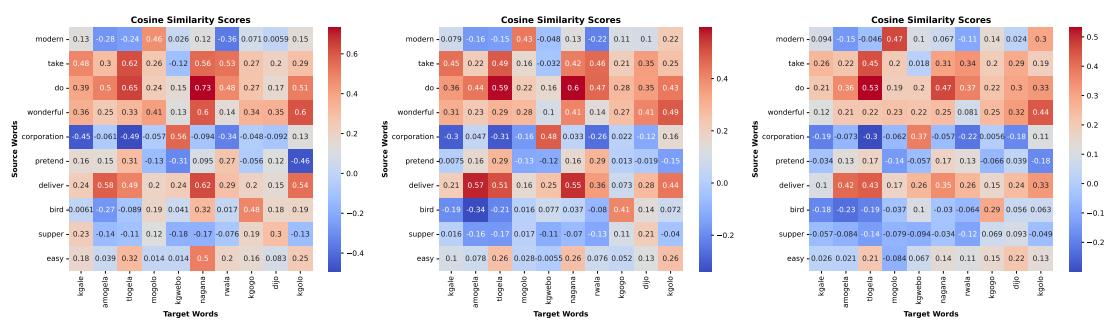
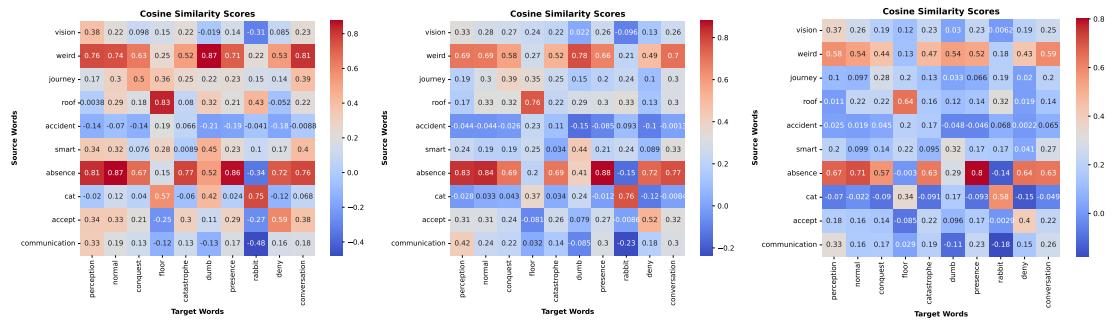


Figure 509: sot VecMap plots of 50 (left), 100(middle), and 200 (right), nso test

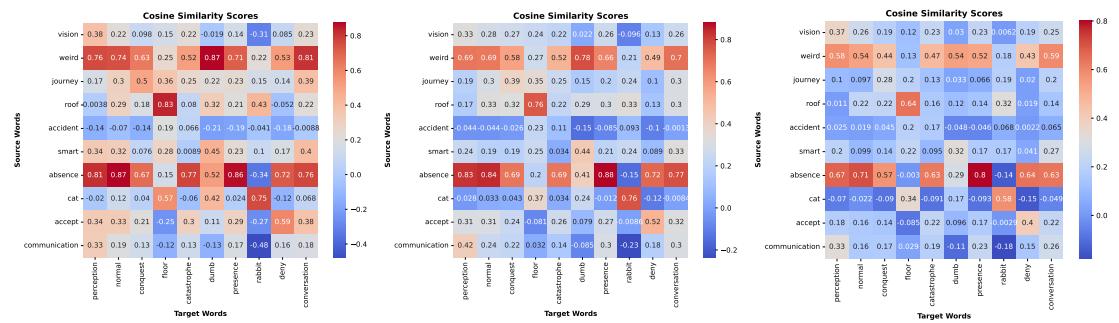


Figure 510: hr

Figure 511: hr

Figure 512: hr

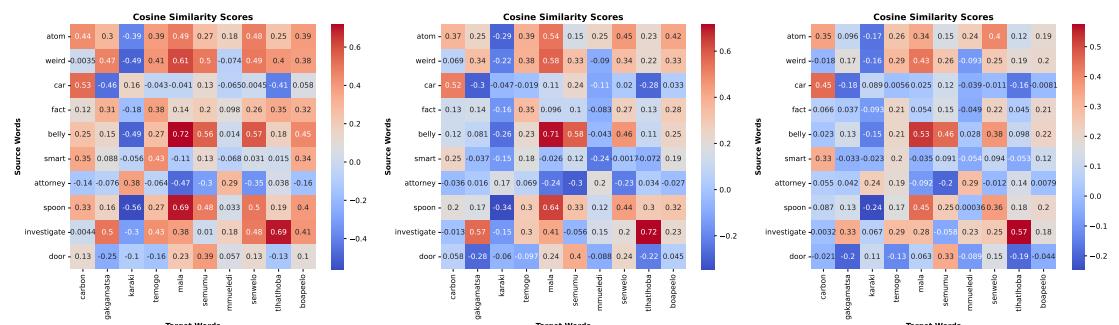


Figure 513: hr and lr-t

Figure 514: hr and lr-t

Figure 515: hr and lr-t

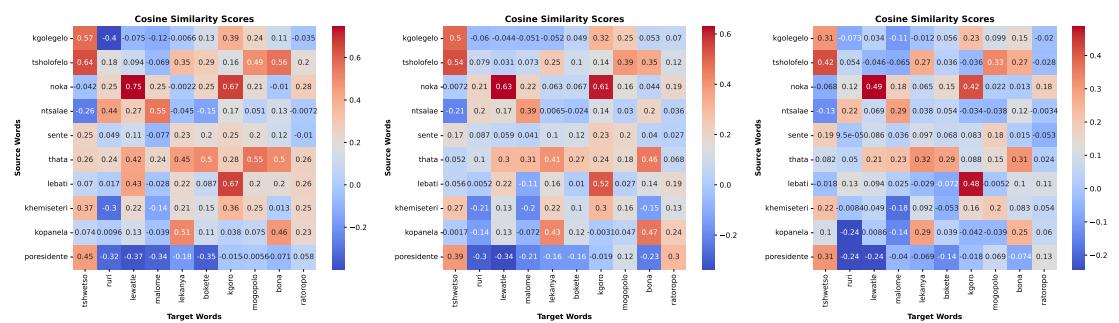


Figure 516: lr-t

Figure 517: lr-t

Figure 518: lr-t

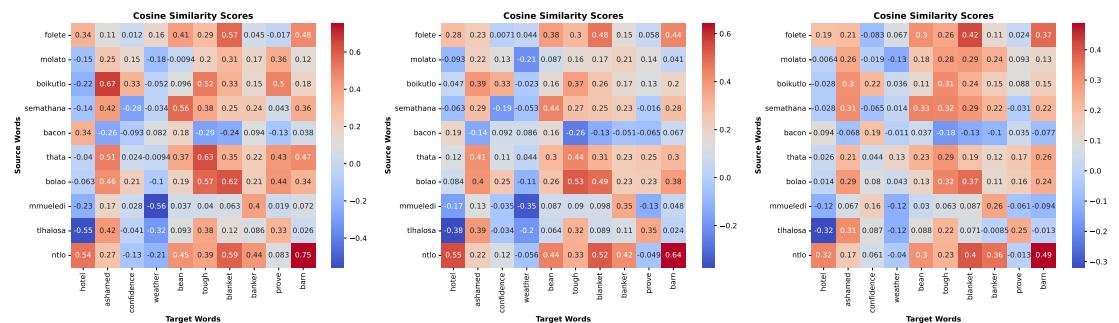


Figure 519: lr-t and hr

Figure 520: lr-t and hr

Figure 521: lr-t and hr

Figure 522: sot VecMap plots of 50 (left), 100(middle), and 200 (right), tsn test

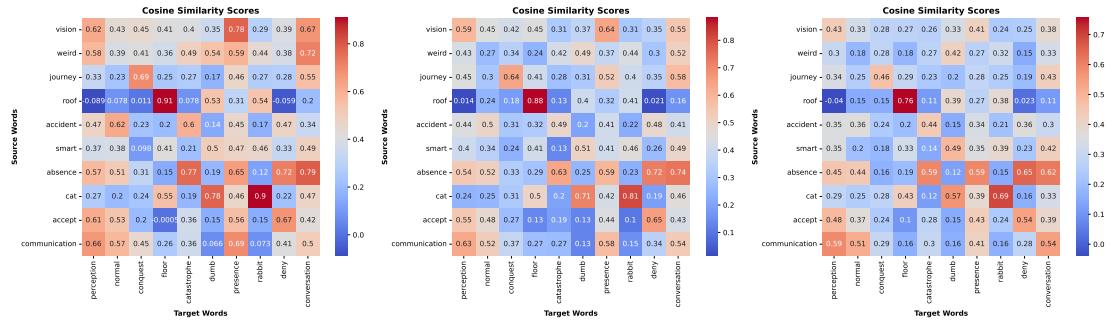


Figure 523: hr

Figure 524: hr

Figure 525: hr

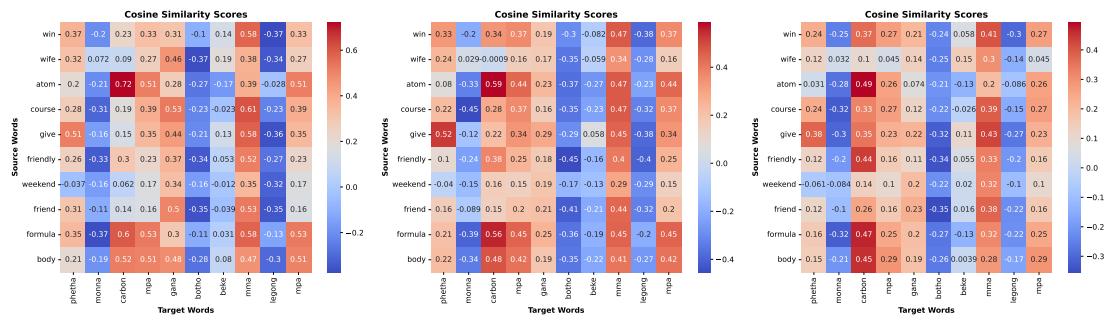


Figure 526: hr and lr-t

Figure 527: hr and lr-t

Figure 528: hr and lr-t

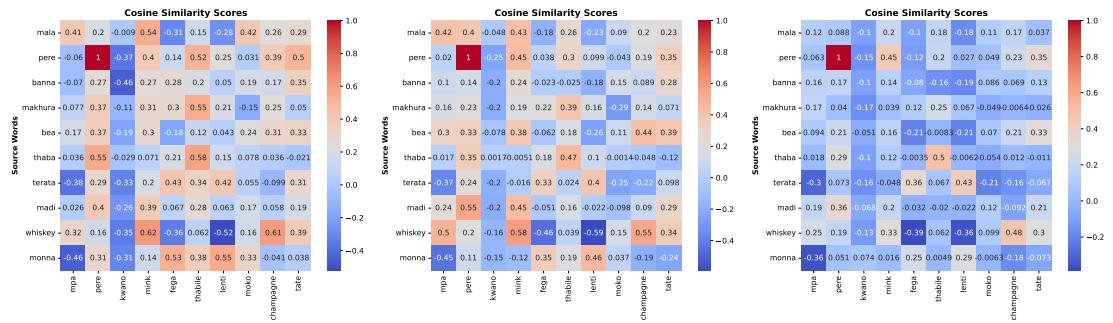


Figure 529: lr-t

Figure 530: lr-t

Figure 531: lr-t

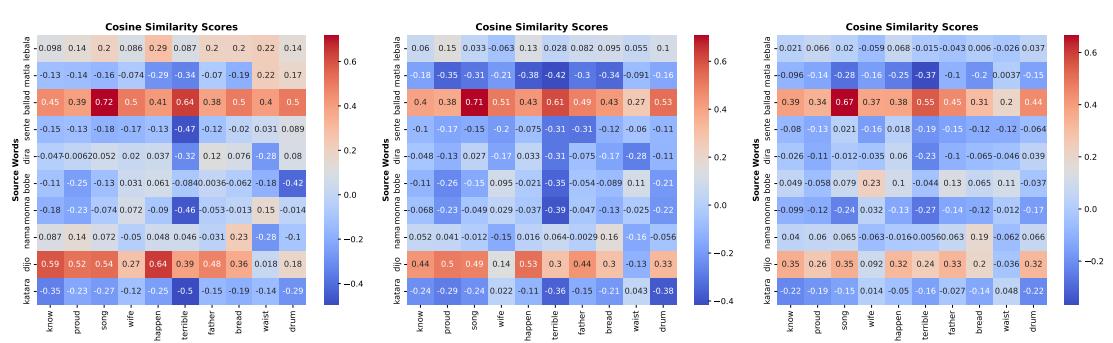


Figure 532: lr-t and hr

Figure 533: lr-t and hr

Figure 534: lr-t and hr

Figure 535: xho VecMap plots of 50 (left), 100(middle), and 200 (right), nso test

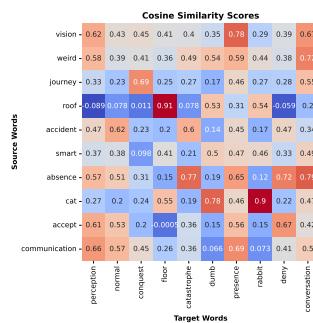


Figure 536: hr

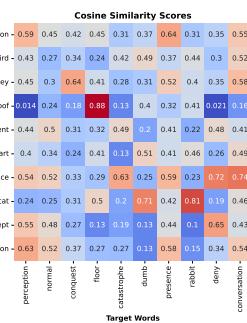


Figure 537: hr

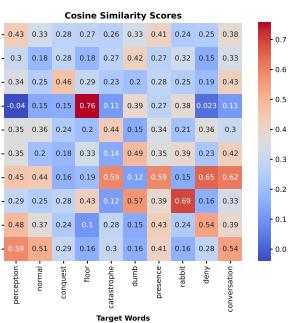


Figure 538: hr

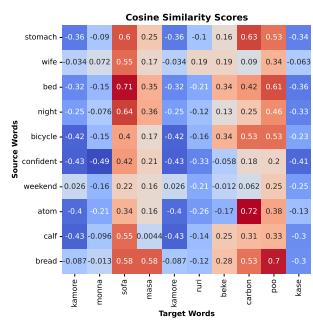


Figure 539: hr and lr-t

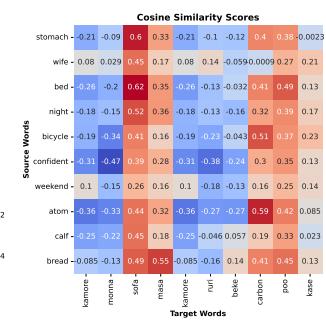


Figure 540: hr and lr-t

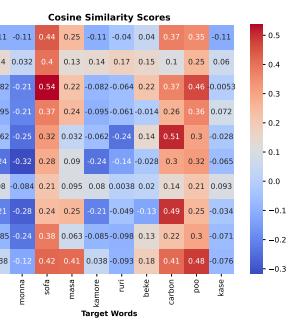


Figure 541: hr and lr-t

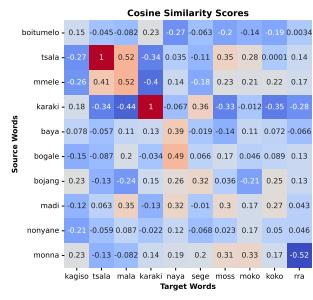


Figure 542: lr-t

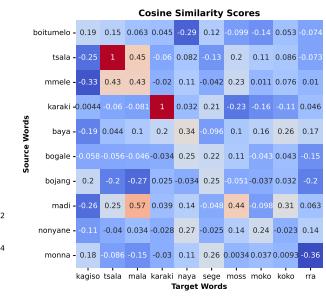


Figure 543: lr-t

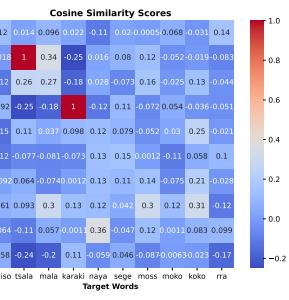


Figure 544: lr-t

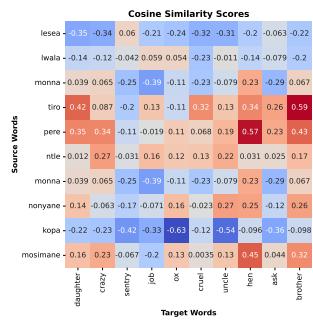


Figure 545: lr-t and hr

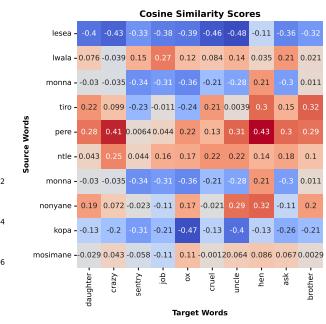


Figure 546: lr-t and hr

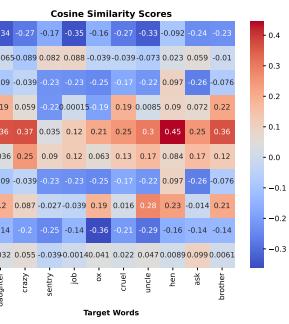


Figure 547: lr-t and hr

Figure 548: xho VecMap plots of 50 (left), 100(middle), and 200 (right), tsn test

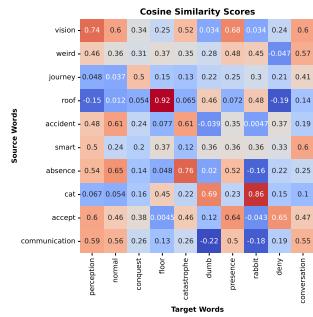


Figure 549: hr

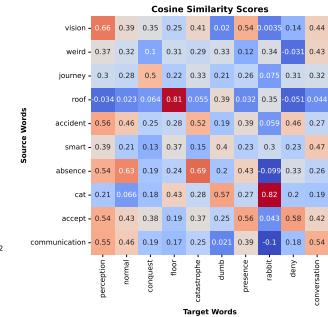


Figure 550: hr

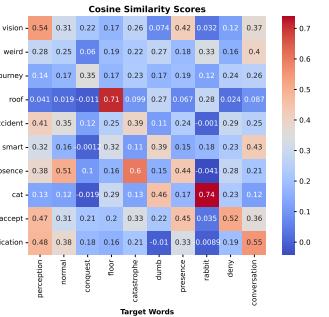


Figure 551: hr

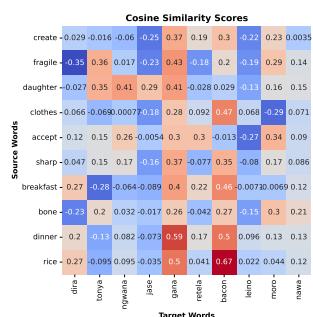


Figure 552: hr and lr-t

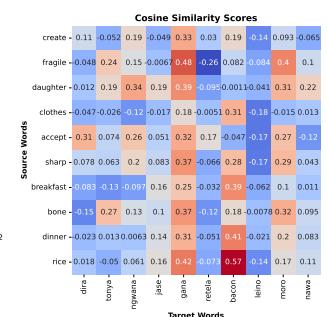


Figure 553: hr and lr-t

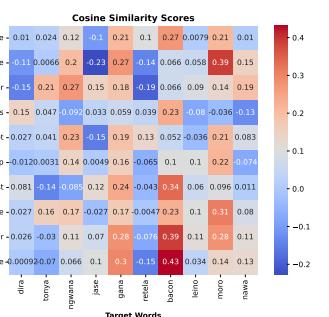


Figure 554: hr and lr-t

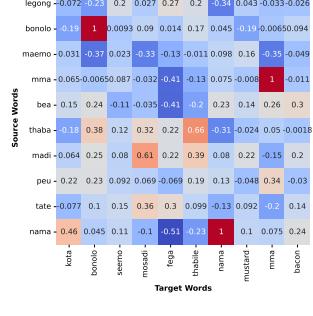


Figure 555: lr-t

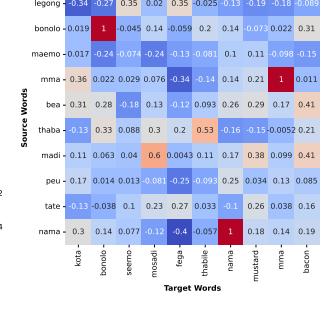
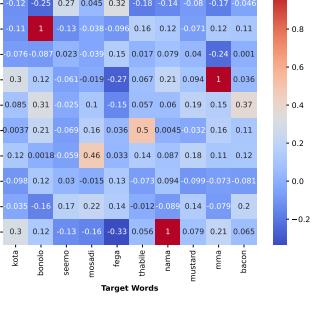


Figure 556: lr-t



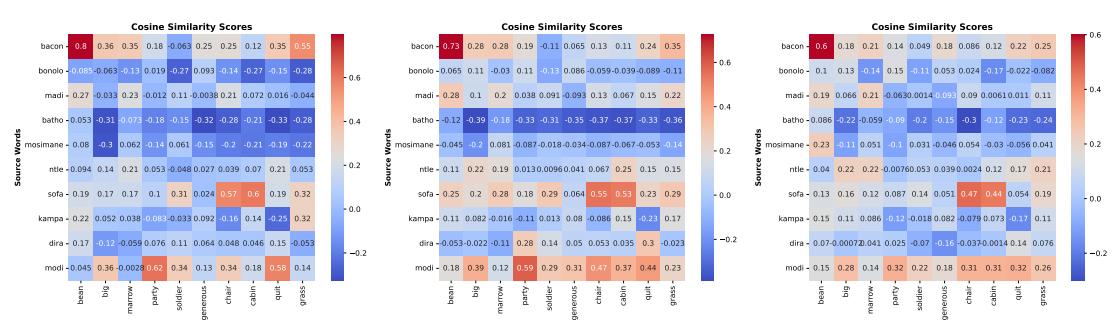
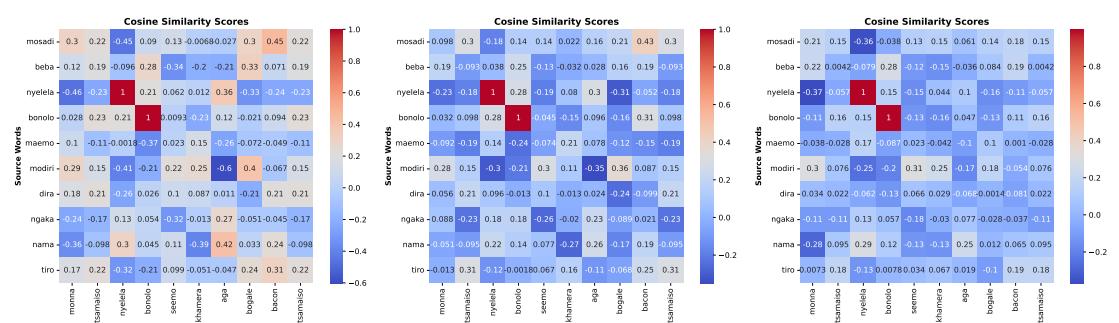
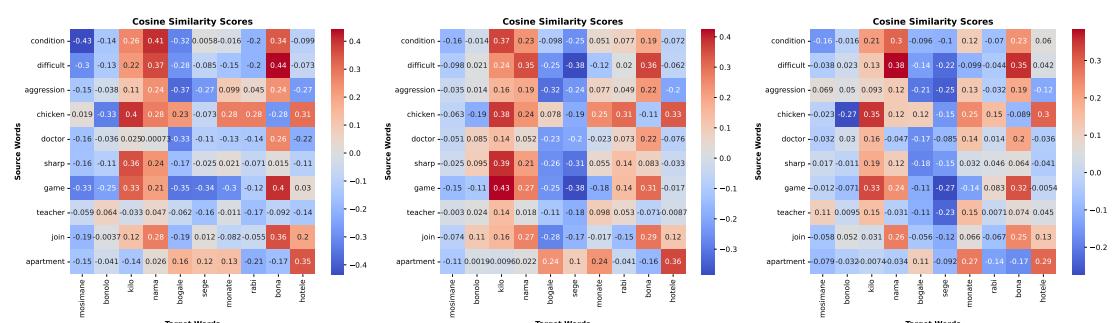
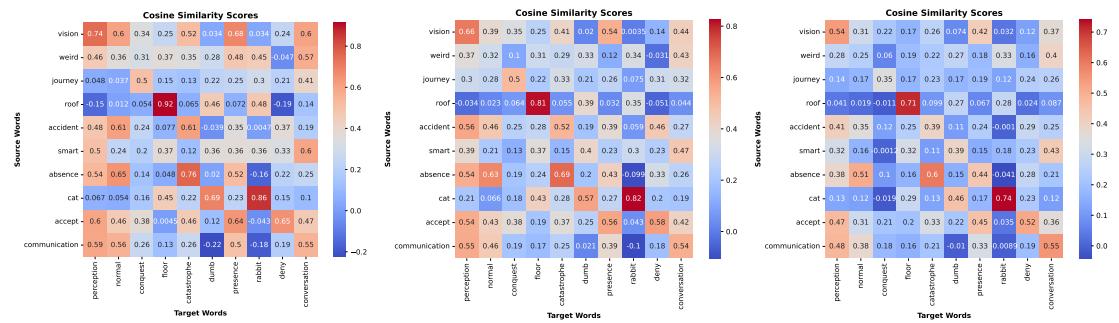


Figure 574: zul VecMap plots of 50 (left), 100(middle), and 200 (right), tsn test

1002 Setswana (tsn) cross-lingual embeddings when us-
1003 ing Muse cross-lingual projections followed by
1004 VecMap embeddings. CCA embeddings generated
1005 the least performing Spearman’s scores. These
1006 results are consistent with previous cosine simi-
1007 larity results outlining the semantic closeness of
1008 embeddings created using the Muse technique com-
1009 pared to CCA and VecMap embeddings. Finally,
1010 although the spearman scores in the plots are low,
1011 they still manage to capture the catapult of mov-
1012 ing from monolingual representation semantics to
1013 cross-lingual representation semantics illustrating
1014 the need to further investigate cross-lingual projec-
1015 tion techniques.

C Cross-lingual and Monolingual Results extended

1016 Table 7, 8, and 9, reports the Accuracy (acc), Preci-
1017 sion (prec), and Recall (rec), of Afro-XLM-r base
1018 model trained on explicit CCA, VecMap, and Muse
1019 cross-lingual embeddings respectively. Our results
1020 indicate that the models perform well under ac-
1021 curacy evaluation, however, they do not do very
1022 well with other metrics such as precision and recall.
1023 This could be due to accuracy being misleading
1024 under imbalanced classes. For Afro-XLM-r model,
1025 Muse embeddings seem to perform better on some
1026 datasets compared to the rest of the techniques.
1027 Additionally, there is no consistent pattern indicating
1028 which embedding technique or dimension performs
1029 consistently better than the rest. Similarly, Table
1030 10, 11, and 12 show a similar trend for the Serengeti
1031 model. Comparatively, Serengeti performed better
1032 than Afro-XLM-r. We hypothesize this could stem
1033 from efficiently learned weights of the Serengeti
1034 model and the size difference since Serengeti is
1035 significantly larger than Afro-XLMr-r. Moreover,
1036 this cold be due to Serengeti being trained on more
1037 African languages compared to Afro-XLM-r, re-
1038 sulting in better transfer of performance gains.
1039

1040 Our vector matrix composing the embedding
1041 layer (i.e. the injected matrix) was created using
1042 randomly selected word vectors from the cross-
1043 lingual embedding word-vector pairs. We hypoth-
1044 esized that this could be improved by selecting
1045 words from the training dataset as word vectors to
1046 build the matrix. For this, Table 13, 14, and 15,
1047 show results of the injected Afro-XLM-r models
1048 created by using words from the code switch train-
1049 ing datasets. Results show no improvements on
1050 the randomly selected word vectors. We tested the
1051

same experiments on monolingual embeddings and
1052 the results are reported in Table 16 with a similar
1053 pattern. We have no clear explanation as to why the
1054 results show no further improvements when using
1055 vectors of the training data as injection embeddings,
1056 and therefore leave this for future works.

D Cross-lingual heatmap scores analyses

1057 In this section, we highlight a few instances where
1058 token attention analyses highlight the consideration
1059 of external words outside the focal sentence for im-
1060 proved text processing. That is, certain words out-
1061 side the focal sentence show high attention scores
1062 compared to words constituting the observed sen-
1063 tence resulting in correct predictions. For exam-
1064 ple, Figure 631, shows the attention heatmap of
1065 the code-switched sentence ‘O utlwile ore dineo o
1066 rileng itll take the pressure off’ with only the high-
1067 est priority score of 0.14 for local words, while the
1068 same sentence shows higher scores of 0.8 when
1069 considering a global embedding space (all embed-
1070 dings) in Figure 632. This means that richer context
1071 is gained externally from the word set of the focal
1072 sentence, arguably made possible by the extended
1073 semantically connected vocabulary derived from
1074 the observed languages. We believe that this is
1075 made possible by the projection of the embedding
1076 space of English and Setswana into the same shared
1077 space, thus extending the context of the overall
1078 embedding space, since the observed pattern does
1079 not exist in the combined but semantically disjoint
1080 monolingual embeddings space of English (Glove)
1081 and Setswana (FastText) embeddings (see Figure
1082 634 and 635 of local attention and global attention
1083 heatmaps respectively).

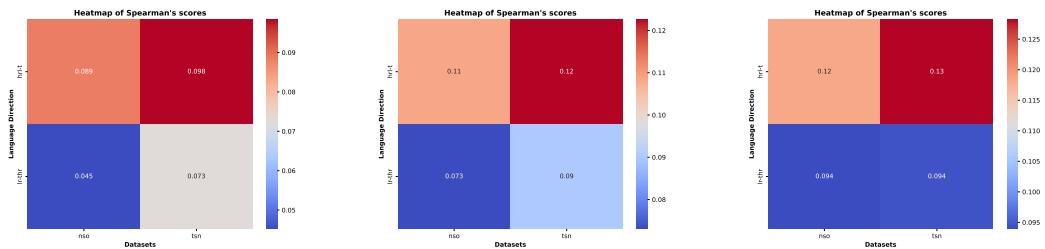


Figure 575: dim=50

Figure 576: dim=100

Figure 577: dim=200

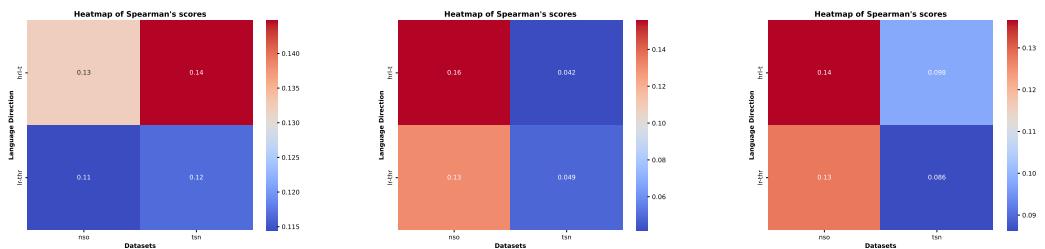


Figure 578: dim=50

Figure 579: dim=100

Figure 580: dim=200

Figure 581: Spearman's correlation: tsn Mono Emb on nso test (row 1) and tsn test (row 2)

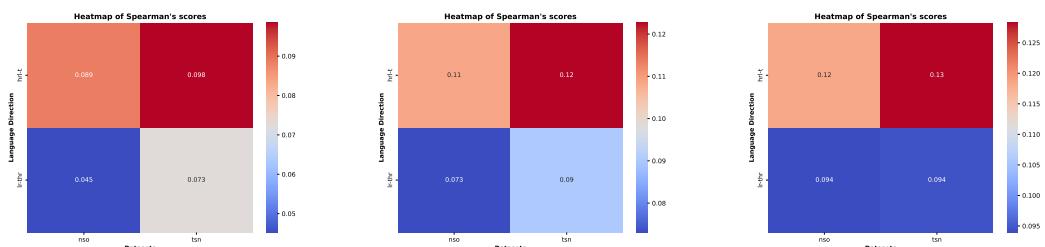


Figure 582: dim=50

Figure 583: dim=100

Figure 584: dim=200

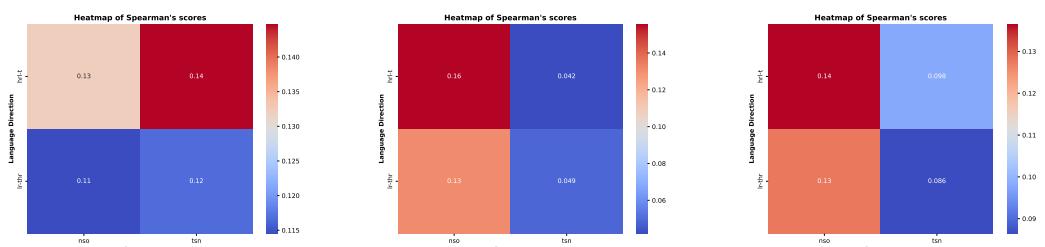


Figure 585: dim=50

Figure 586: dim=100

Figure 587: dim=200

Figure 588: Spearman's correlation: sot Mono Emb on nso test (row 1) and tsn test (row 2)

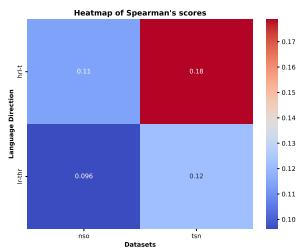


Figure 589: dim=50

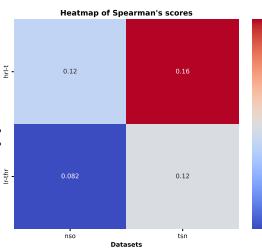


Figure 590: dim=100

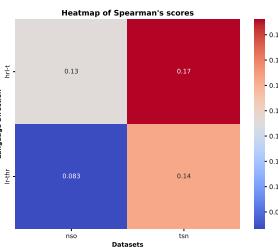


Figure 591: dim=200

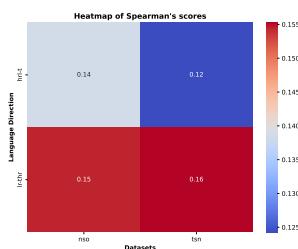


Figure 592: dim=50

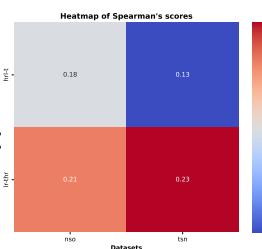


Figure 593: dim=100

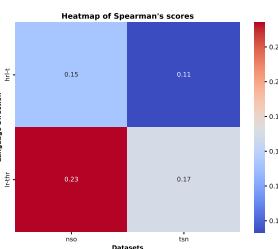


Figure 594: dim=200

Figure 595: Spearman's correlation: en-tsn CCA cross Emb on nso test (row 1) and tsn test (row 2)

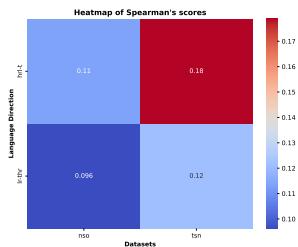


Figure 596: dim=50

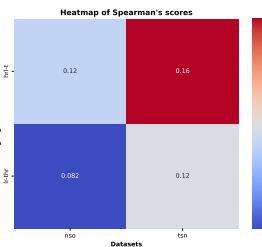


Figure 597: dim=100

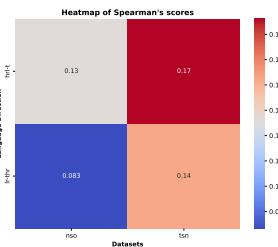


Figure 598: dim=200

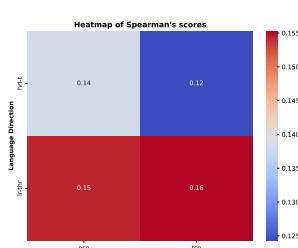


Figure 599: dim=50

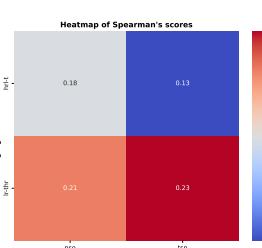


Figure 600: dim=100

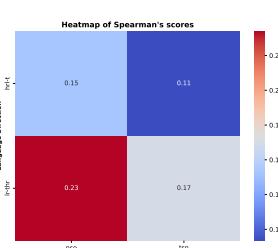


Figure 601: dim=200

Figure 602: Spearman's correlation: en-sot CCA cross Emb on nso test (row 1) and tsn test (row 2)

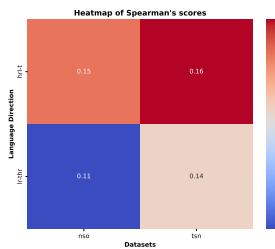


Figure 603: dim=50

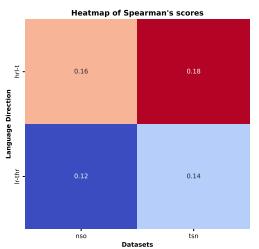


Figure 604: dim=100

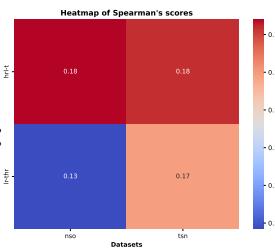


Figure 605: dim=200

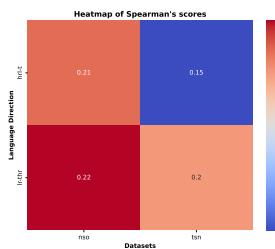


Figure 606: dim=50

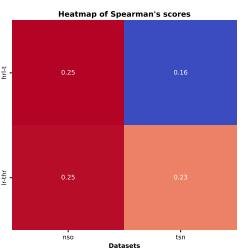


Figure 607: dim=100

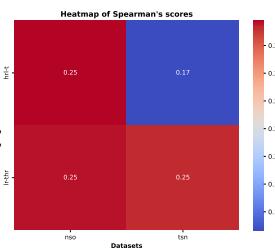


Figure 608: dim=200

Figure 609: Spearmans correlation: en-tsn Muse cross Emb on nso test (row 1) and tsn test (row 2)

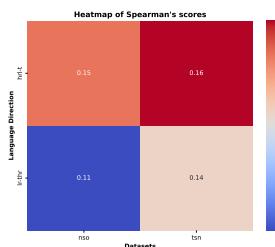


Figure 610: dim=50

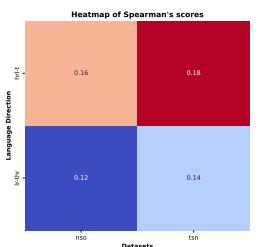


Figure 611: dim=100

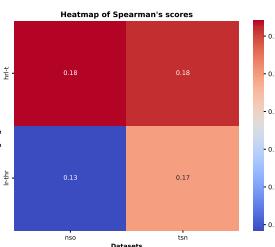


Figure 612: dim=200

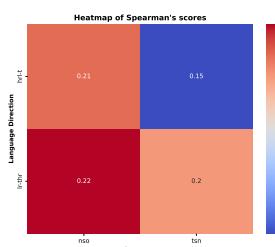


Figure 613: dim=50

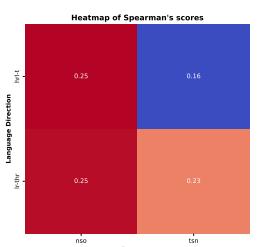


Figure 614: dim=100

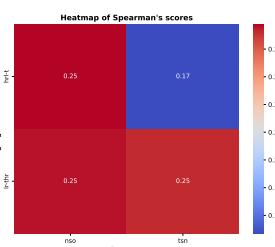


Figure 615: dim=200

Figure 616: Spearmans correlation: en-sot Muse cross Emb on nso test (row 1) and tsn test (row 2)

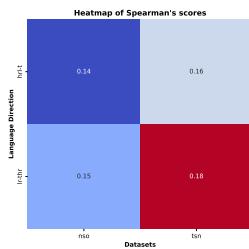


Figure 617: dim=50

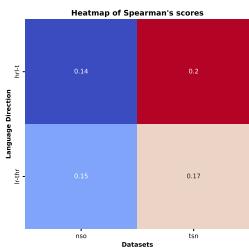


Figure 618: dim=100

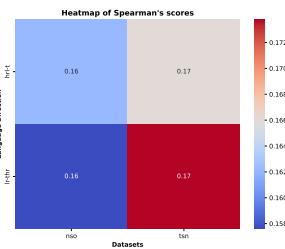


Figure 619: dim=200

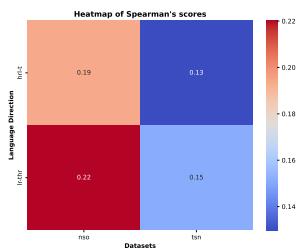


Figure 620: dim=50

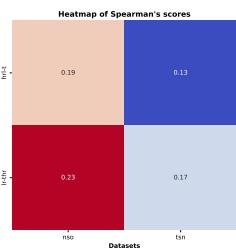


Figure 621: dim=100

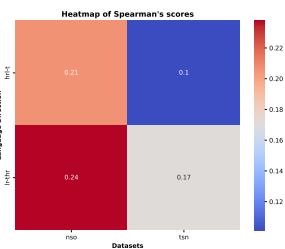


Figure 622: dim=200

Figure 623: Spearmans correlation: en-tsn VecMap cross Emb on nso test (row 1) and tsn test (row 2)

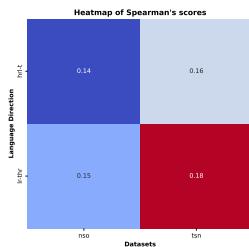


Figure 624: dim=50

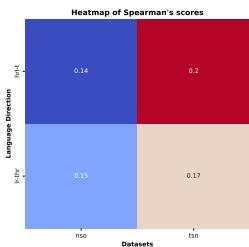


Figure 625: dim=100

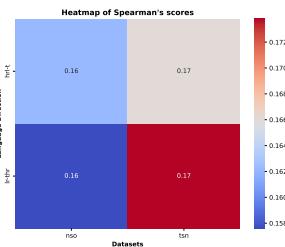


Figure 626: dim=200

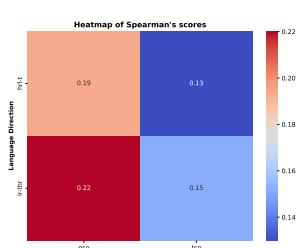


Figure 627: dim=50

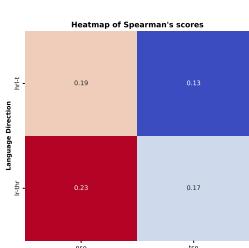


Figure 628: dim=100

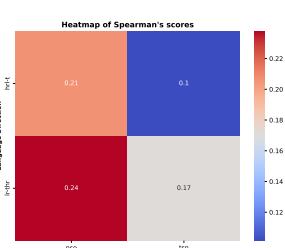


Figure 629: dim=200

Figure 630: Spearmans correlation: en-sot VecMap cross Emb on nso test (row 1) and tsn test (row 2)

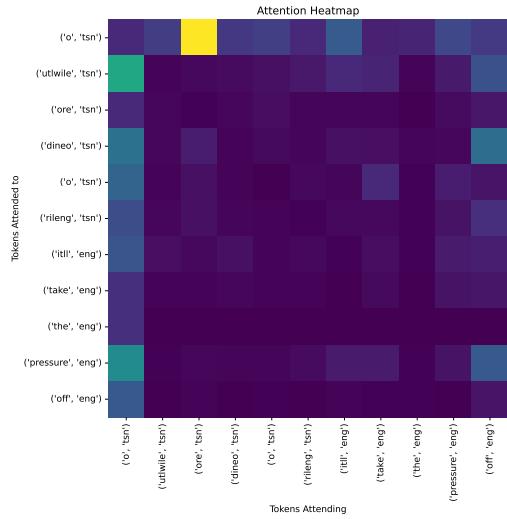


Figure 631: VecMap local attention

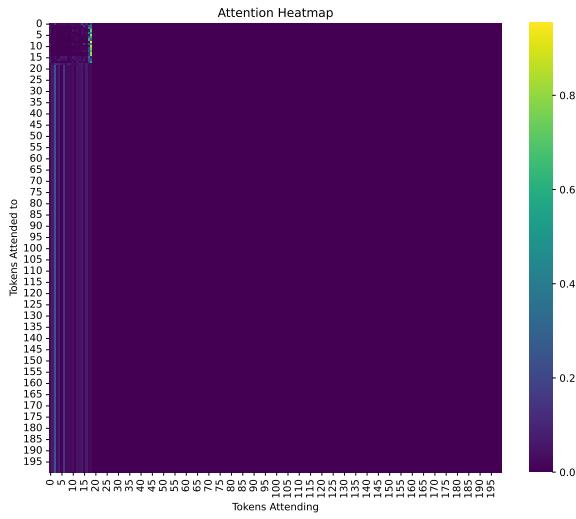


Figure 632: VecMap global attention

Figure 633: Embedding Attention heatmap visualization of Cross-lingual (CL) von Setswana (tsn) sentence

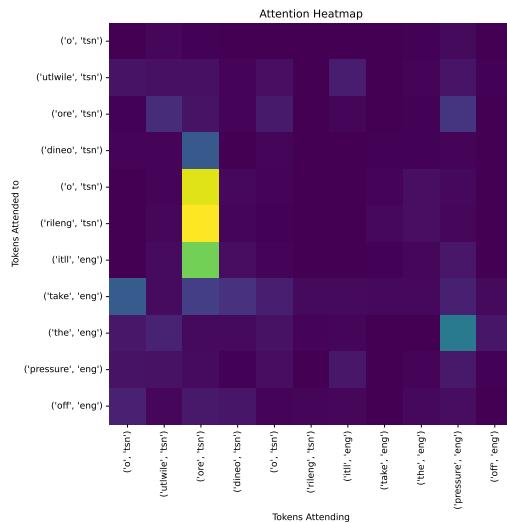


Figure 634: Mono local attention

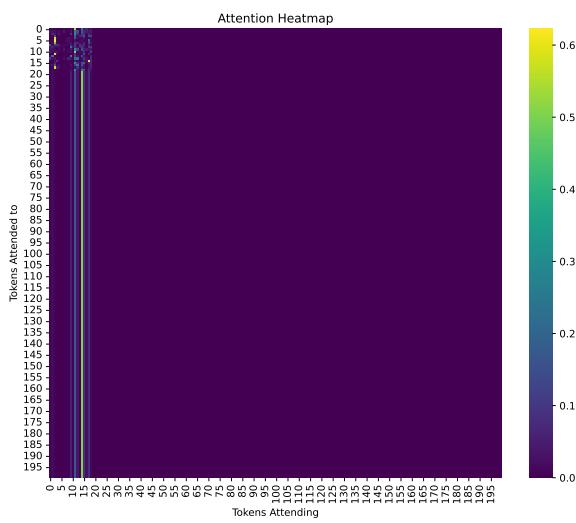


Figure 635: Mono global attention

Figure 636: Embedding Attention heatmap visualization of Cross-lingual (CL) von Setswana (tsn) sentence

	sot	tsn	zul	xho	zulxho	sottsn	zulxhosot	zulxhotsn	sottsnzul	sottsnxho	zulxhosottsn
Embedding Dimension: 50											
acc	97	97	94	94	76	74	90	81	80	89	89
prec	81	84	80	82	59	61	75	65	67	77	78
rec	82	85	75	78	56	55	73	59	66	79	80
Embedding Dimension: 100											
acc	96	97	89	82	85	86	91	92	89	89	72
prec	80	83	79	69	70	73	76	78	79	77	59
rec	82	84	75	64	69	76	73	74	81	79	55
Embedding Dimension: 200											
acc	86	86	94	88	84	81	69	85	89	87	77
prec	73	65	80	75	69	68	53	72	79	73	63
rec	71	59	76	69	68	67	46	71	81	71	58

Table 7: Afro-XLMr-base + CCA : Average metrics for each embedding dimension across folders. Averaged over 5 runs. The word embeddings for injection created through random selection of words to build the vector matrix.

	sot	tsn	zul	xho	zulxho	sottsn	zulxhosot	zulxhotsn	sottsnzul	sottsnxho	zulxhosottsn
Embedding Dimension: 50											
acc	95	97	97	97	93	78	95	93	80	79	63
prec	76	83	85	85	80	65	83	78	68	66	49
rec	75	85	86	83	81	65	84	75	67	61	43
Embedding Dimension: 100											
acc	97	97	95	97	85	86	95	68	88	81	79
prec	81	83	79	85	70	73	82	49	75	69	63
rec	82	85	77	84	69	76	83	44	76	68	62
Embedding Dimension: 200											
acc	86	97	96	90	90	86	94	85	89	77	89
prec	72	83	80	78	73	74	79	71	77	63	77
rec	71	85	79	74	72	77	80	69	79	57	79

Table 8: Afro-XLMr-base + VecMap: Average metrics for each embedding dimension across folders. Averaged over 5 runs. The word embeddings for injection are created through a random selection of words to build the vector matrix.

	sot	tsn	zul	xho	zulxho	sottsn	zulxhosot	zulxhotsn	sottsnzul	sottsnxho	zulxhosottsn
Embedding Dimension: 50											
acc	96	97	86	88	90	85	92	93	78	79	87
prec	75	79	65	69	69	69	71	72	59	62	67
rec	78	80	62	67	72	72	75	76	61	63	72
Embedding Dimension: 100											
acc	85	97	95	94	91	85	93	84	87	88	87
prec	68	80	75	74	70	69	72	64	68	69	67
rec	67	82	77	75	74	73	76	65	73	75	73
Embedding Dimension: 200											
acc	96	97	95	95	90	86	93	93	88	88	88
prec	74	79	75	74	70	69	72	73	69	69	69
rec	78	81	78	76	73	74	77	78	75	74	75

Table 9: Afro-XLMr-base + Muse: Average metrics for each embedding dimension across folders. Averaged over 5 runs. The word embeddings for injection are created through a random selection of words to build the vector matrix.

	sot	tsn	zul	xho	zulxho	sottsn	zulxhosot	zulxhotsn	sottsnzul	sottsnxho	zulxhosottsn
Embedding Dimension: 50											
acc	96	97	96	96	93	86	94	95	89	89	89
prec	77	78	79	80	78	72	79	80	76	76	76
rec	81	82	82	82	80	76	81	82	79	78	79
f1	79	80	80	81	79	74	80	81	78	77	78
Embedding Dimension: 100											
acc	96	97	96	96	93	87	94	95	89	89	89
prec	77	79	78	81	78	72	80	80	76	76	77
rec	81	82	81	82	79	76	82	82	79	79	79
f1	78	80	80	82	79	74	81	81	78	77	78
Embedding Dimension: 200											
acc	96	97	96	96	93	87	94	95	89	89	90
prec	77	80	79	81	78	73	80	81	76	77	77
rec	81	83	81	82	80	77	82	83	79	80	79
f1	79	82	80	82	79	75	81	82	78	78	78

Table 10: Serengeti + CCA: Average metrics for each embedding dimension across folders. Averaged over 5 runs. The word embeddings for injection are created through random selection of words to build the vector matrix.

	sot	tsn	zul	xho	zulxho	sottsn	zulxhosot	zulxhotsn	sottsnzul	sottsnxho	zulxhosottsn
Embedding Dimension: 50											
acc	96	97	96	-	93	86	94	95	89	89	89
f1	78	80	79	-	78	74	80	81	77	77	77
prec	76	79	78	-	77	72	78	80	76	76	76
rec	81	82	81	-	79	76	81	82	79	79	79
Embedding Dimension: 100											
acc	96	97	96	-	93	87	94	95	89	89	89
prec	77	79	79	-	78	72	79	81	76	76	76
rec	81	83	81	-	79	76	81	82	79	79	79
f1	79	81	80	-	79	74	80	81	77	77	78
Embedding Dimension: 200											
acc	96	97	96	-	93	87	94	95	89	89	90
prec	77	80	79	-	79	73	80	82	77	76	78
rec	81	82	81	-	80	77	82	83	79	79	80
f1	79	81	80	-	79	75	81	82	78	78	79

Table 11: Serengeti + VecMap: Average metrics for each embedding dimension across folders. Averaged over 5 runs. The word embeddings for injection are created through a random selection of words to build the vector matrix.

	sot	tsn	zul	xho	zulxho	sottsn	zulxhosot	zulxhotsn	sottsnzul	sottsnxho	zulxhosottsn
Embedding Dimension: 50											
acc	96	96	95	95	91	86	93	94	88	88	88
prec	73	75	74	76	71	68	73	73	70	69	70
rec	79	81	78	79	75	75	78	79	75	76	75
f1	76	78	76	77	73	72	75	76	72	73	72
Embedding Dimension: 100											
acc	96	96	95	96	92	86	93	94	88	88	88
prec	73	75	74	77	73	67	73	74	69	70	70
rec	80	81	79	80	77	75	79	79	76	76	76
f1	76	78	77	78	75	71	76	77	73	73	73
Embedding Dimension: 200											
acc	96	96	95	96	92	86	93	94	88	88	88
prec	73	74	75	77	73	68	74	75	70	70	71
rec	80	81	78	80	77	75	79	80	76	76	77
f1	76	77	76	78	75	71	77	77	73	73	74

Table 12: Serengeti + Muse: Average metrics for each embedding dimension across folders. Averaged over 5 runs. The word embeddings for injection are created through a random selection of words to build the vector matrix.

	sot	tsn	zul	xho	zulxho	sottsn	zulxhosot	zulxhotsn	sottsnzul	sottsnxho	zulxhosottsn
Embedding Dimension: 50											
acc	94	82	84	94	76	72	67	76	73	89	65
prec	76	61	66	80	58	59	51	60	58	75	51
rec	72	57	55	76	55	53	45	57	54	77	44
f1	74	58	60	78	56	55	47	58	55	76	46
Embedding Dimension: 100											
acc	94	79	81	84	67	69	62	85	71	69	72
prec	75	55	71	70	47	56	45	71	56	57	59
rec	72	45	63	60	43	52	35	69	53	48	55
f1	73	48	66	64	44	54	39	70	54	51	56
Embedding Dimension: 200											
acc	90	74	89	94	84	60	82	66	71	62	68
prec	71	50	73	80	66	49	65	49	57	50	52
rec	63	43	64	77	64	41	60	43	54	31	44
f1	66	45	67	79	65	44	62	45	55	38	47

Table 13: Afro-XLMR-base + CCA: Average metrics for each embedding dimension across folders. Averaged over 5 runs. The word embeddings for injection are created through word selection from training data to build the vector matrix.

	sot	tsn	zul	xho	zulxho	sottsn	zulxhosot	zulxhotsn	sottsnzul	sottsnxho	zulxhosottsn
Embedding Dimension: 50											
acc	86	85	89	79	67	66	72	56	71	77	63
prec	70	64	76	65	47	53	56	37	56	65	48
rec	68	58	72	58	42	44	48	29	53	58	42
f1	69	60	73	61	44	47	51	32	54	61	45
Embedding Dimension: 100											
acc	83	97	88	82	59	64	71	80	57	74	68
prec	63	82	74	65	36	52	54	63	40	61	52
rec	58	84	70	61	29	43	47	58	29	50	44
f1	60	83	71	63	32	46	49	60	33	55	47
Embedding Dimension: 200											
acc	93	74	81	88	75	75	66	85	65	78	69
prec	70	50	71	74	56	59	50	69	50	64	50
rec	65	42	62	68	53	55	44	68	46	58	42
f1	67	45	65	71	54	57	46	68	47	60	45

Table 14: Afro-XLMR-base + VecMap: Average metrics for each embedding dimension across folders. Averaged over 5 runs. The word embeddings for injection are created through word selection from training data to build the vector matrix.

	sot	tsn	zul	xho	zulxho	sottsn	zulxhosot	zulxhotsn	sottsnzul	sottsnxho	zulxhosottsn
Embedding Dimension: 50											
acc	88	86	91	82	88	85	88	76	78	79	79
prec	63	61	64	62	63	65	65	54	58	59	60
rec	49	55	61	57	66	70	65	52	60	61	61
f1	54	57	62	60	65	68	65	52	58	60	60
Embedding Dimension: 100											
acc	80	73	91	81	81	85	83	91	61	84	78
prec	59	47	67	62	57	68	61	67	42	64	58
rec	46	39	64	58	58	72	51	69	38	64	59
f1	51	42	65	59	57	70	54	68	39	64	58
Embedding Dimension: 200											
acc	82	85	73	79	90	68	73	83	61	79	66
prec	59	57	60	60	69	54	53	62	42	62	45
rec	51	52	48	55	73	51	42	63	38	63	38
f1	54	54	52	57	71	51	46	62	39	62	41

Table 15: Afro-XLMR-base + Muse: Average metrics for each embedding dimension across folders. Averaged over 5 runs. The word embeddings for injection are created through word selection from training data to build the vector matrix.

	sot	tsn	zul	xho	zulxho	sottsn	zulxhosot	zulxhotsn	sottsnzul	sottsnxho	zulxhosottsn
Embedding Dimension: 100											
acc	93	87	72	79	81	75	81	74	78	76	75
prec	63	58	57	61	58	55	55	52	59	59	55
rec	64	56	44	54	60	55	49	51	61	55	52
f1	63	56	49	57	59	55	52	51	60	56	52
Embedding Dimension: 200											
acc	92	81	76	91	72	82	65	65	55	81	75
prec	66	56	56	72	46	63	45	44	37	61	56
rec	64	52	33	66	44	62	41	40	28	63	54
f1	65	53	41	69	45	61	42	41	31	62	54
Embedding Dimension: 50											
acc	93	81	91	91	80	82	87	83	67	76	79
prec	71	56	61	69	56	64	63	61	47	57	59
rec	68	53	62	65	56	62	64	61	40	54	61
f1	69	54	61	67	56	62	64	61	43	55	60

Table 16: Afro-XLMr-base + Monolingual: Average metrics for each embedding dimension across folders. Averaged over 5 runs. The word embeddings for injection are created through a random selection of words to build the vector matrix.