

Zero-shot Cross-lingual Transfer Learning with Multiple Source and Target Languages for Information Extraction: Language Selection and Adversarial Training

Nghia Trung Ngo and Thien Huu Nguyen

Department of Computer Science

University of Oregon, Eugene, Oregon, USA

nghian@uoregon.edu, thien@cs.uoregon.edu

Abstract

The majority of previous researches addressing multi-lingual IE are limited to zero-shot cross-lingual single-transfer (one-to-one) setting, with high-resource languages predominantly as source training data. As a result, these works provide little understanding and benefit for the realistic goal of developing a multi-lingual IE system that can generalize to as many languages as possible. Our study aims to fill this gap by providing a detailed analysis on Cross-Lingual Multi-Transferability (many-to-many transfer learning), for the recent IE corpora that cover a diverse set of languages. Specifically, we first determine the correlation between single-transfer performance and a wide range of linguistic-based distances. From the obtained insights, a combined language distance metric can be developed that is not only highly correlated but also robust across different tasks and model scales. Next, we investigate the more general zero-shot multi-lingual transfer settings where multiple languages are involved in the training and evaluation processes. Language clustering based on the newly defined distance can provide directions for achieving the optimal cost-performance trade-off in data (languages) selection problem. Finally, a relational-transfer setting is proposed to further incorporate multi-lingual unlabeled data based on adversarial training using the relation induced from the above linguistic distance. Experimental results on two practical multi-lingual IE tasks demonstrate our method significantly outperforms baselines across tasks and languages simultaneously. Additionally, by carefully designing the multi-lingual training to utilize data from relevant languages, we can achieve a substantial boost in generalization ability with reasonable labor cost for the additional data collection.

1 Introduction

The objective of Information extraction (IE) is to identify and extract structure information, such

as entities, relations, and events, from natural unstructured text. IE plays an important role in various downstream applications, including Question Answering, Knowledge Graph Construction, News Analysis, etc. Solving IE tasks pose significant challenges for NLP models as they often require the understanding of complex features of natural languages. For example, to extract relations within a sentence, models first need to learn specialized structures of the corresponding language to identify entities mentioned in the given text. Next, a deep understanding of context is required to correctly classify the relations between these entities. These challenges are further exacerbated in multilingual settings, where datasets are collected from multiple languages, each of which contains language-specific characteristics and structures.

The rapid development around English-based datasets has pushed machine performance to be on par with human ability in English tasks, prompting recent works to explore NLP research in other languages (Liang et al., 2020; Ruder et al., 2021). However, despite advanced large-scale architectures and high English results, current models notably under-perform in new languages, especially those that are considered low-resource and lack high-quality datasets for ne-tuning. Cross-lingual Transfer, as a result, becomes one of the most important directions in the field. Given a particular task, the goal of Cross-lingual Transfer is to train multilingual models over high-resource source languages that can solve textual tasks in new target languages despite the shifts in linguistic origin.

Currently, the most popular and practical approach for IE involves Zero-Shot Cross-Lingual (ZSCL) transfer (Conneau et al., 2020; Goyal et al., 2021). These methods fine-tune Transformer-based multilingual Language Models (mLMs), which were pre-trained using unlabeled text from hundreds of languages, for downstream tasks using high-resource source-language labeled

datasets (predominantly English). The resulting models are directly used for evaluation on the corresponding tasks in target languages. Studies have shown, however, performance of these multi-lingual models varies substantially across languages and tasks. Several factors have been attributed to this phenomenon, ranging from data-dependent statistic (e.g. dataset size, word overlap) (Malkin et al., 2022), to data-independent features (e.g. phylogenetic and typological features) (Lin et al., 2019; Dolicki and Spanakis, 2021). Based on these previous observations, we believe that there is a deeper connection between cross-lingual transfer ability and the relations in the linguistic landscape. Unraveling this correlation can provide tremendous practical implications for IE. First, it serves as a guide for data collection process to achieve optimal cost-performance trade-off by gathering training samples from appropriate source languages for a target language. Furthermore, the modeling process can also be tailored such that the learned representations explicitly capture the linguistic relations to improve generalization across languages.

Previous papers following the above direction define the problem as a Performance Prediction task. In (Lin et al., 2019; Dolicki and Spanakis, 2021; Srinivasan et al., 2021), a regression model is trained to take linguistic features of the source-target language pairs as input to predict a trained model performance scores on target languages. Despite high prediction accuracy, these works are insufficient for the following two reasons. First, they place too much emphasis on the accuracy of the regression model, which is trained for a specific architecture on a particular task. As the training configuration varies widely in practice, the results obtained from the performance prediction models may become unreliable and not applicable in general. Another reason is that previous work is only limited to the setting of single-transfer between two languages, in which only one source language (predominately English) is utilized. Current advances in translation model and data gathering process have enabled the creation of datasets in many languages, thus multiple source languages should be considered. Intuitively, additional training data from more languages can help improve model’s generalization on downstream tasks, and learning from text in multiple languages may have a positive effect on zero-shot transfer. We be-

lieve that multi-transfer setting is the next important step for cross-lingual transfer, both to improve model performance across languages and to provide a more complete picture of multi-linguality in machine learning. In this paper, we focus on what has been missing in previous works by aiming to answer the following three major research questions:

Q1: How do the relations in the linguistic landscape affect an IE models cross-lingual transfer ability? We use URIEL Typological Database (Littell et al., 2017) to extract phylogenetic and typological properties of each language. These properties, represented as multi-dimensional binary features, are used to compute the pair-wise linguistic distances (or equivalently the similarity scores) between languages. We compare the correlation between these scores and model single-transfer performance. A source language with a high correlation value would imply that we can infer its ability to transfer to different languages using only linguistic relations, without the need to actually fine-tune models.

Q2: Can we implicitly leverage these linguistic features as dataset-independent knowledge to efficiently address the more general multi-transfer setting? While many-to-many cross-lingual transfer has the potential to significantly improve single-transfer performance, it would also require gathering data from multiple languages. Given a set of languages and their linguistic features as the only prior information, we aim to find the optimal subset of source languages to gather labeled data for zero-shot cross-lingual multi-transfer to target languages. The goal is to achieve the best cost-performance trade-off on all languages, without having to fine-tune the mLMs on an exponential number of possible language combinations.

Q3: Can we explicitly integrate these linguistic relations in the learning process to effectively improve multi-transfer performance? We then investigate further into the possibility of directly embedding the linguistic relations in the fine-tuning process. The hypothesis is that, by capturing these connections, the multi-lingual representations would be able to adaptively generalize to not only languages that are closely related to source languages, but also distant languages that share little similarity with the available training data.

The following observations are obtained from our quantitative experiments and qualitative analysis, through 3 levels of transfer settings:

1) Single-transfer (ZSCL-S) - Only 1 labeled source language available. It is possible to achieve a high degree of correlation between model ZSCL performances and linguistic relations, using a combination of syntax, inventory, and phonology features from URIEL. However, in contrast to prior works which only focus on syntactic transfer when fine-tuning, our combined metric places the least importance on the syntax feature. This implies that previous researches are suboptimal and incomplete, prompting further investigations into the problem.

2) Multi-transfer (ZSCL-M) - Multiple labeled source language available. We first cluster languages based on the combined metric above. Then, by selecting source languages following the guidance from the resulting clusters, we observe significant improvements in ZSML performances over the naive method of randomly picking source languages. In other words, with only the prior linguistic knowledge, we can efficiently choose a suitable small subset of languages for labeled data annotations, to fine-tune a MMLM to perform best on a given set of languages.

3) Relational-transfer (ZSCL-R) - ZSCL-M with additional multi-lingual unlabeled data. We leverage unlabeled data from all available languages and their linguistic relations as inputs to graph-relational adversarial learning framework (Xu et al., 2022b), a generalization of adversarial language adaption (Chen et al., 2018) that can only perform strict uniform alignment for pair-wise transfer. By conditioning the multi-lingual representation flexibly on the connections expressed by the corresponding language relational-graph, we achieve a considerable increase in transfer performances across every language. This is only at the small cost of collecting additional unlabeled data from other languages.

2 Related Work

Zero-shot Cross-lingual Transfer The majority of recent ZSCL works (Fang et al., 2021; Chi et al., 2021) follow single-transfer setting, using comprehensive multi-lingual multi-task benchmarks such as XTREME (Ruder et al., 2021), or XGLUE (Liang et al., 2020). These datasets provide English training data for fine-tuning the pre-trained

SMILER		MINION	
	%		%
ita	19.71	eng	39.76
fra	16.25	pol	13.7
deu	13.75	tur	13.7
por	11.54		
nld	10.38		
eng	9.57	spa	9.99
kor	5	por	4.59
pol	4.5	swe	4.59
spa	2.95		
ara	2.49		
rus	1.71	hin	4.58
swe	1.2	kor	4.58
fas	0.7	jpn	4.5
ukr	0.26		

Table 1: Percentage distributions of training data in each task for every language, which are separated into high, medium, and low resource categories. The shared languages are color-coded, with red indicating that the language belongs to a different category between the two tasks, whereas green indicates otherwise. This study involves a total of 17 languages including: arabic (ara), german (deu), english (eng), farsia (fas), french (fra), hindi (hin), italian (ita), japanese (jpn), korean (kor), dutch (nld), polish (pol), portuguese (por), russian (rus), spanish (spa), swedish (swe), turkish (tur), and ukrainian (ukr).

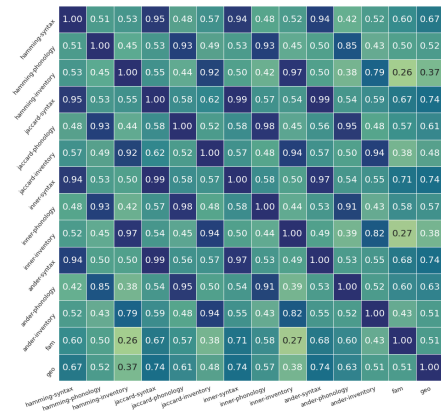


Figure 1: The pairwise Pearson correlation for all computed language distances.

MMLMs, which are then evaluated on translated test sets in different languages. As a result, English becomes the dominant source language for transfer in following ZSCL researches (Phang et al., 2020), which is suboptimal due to the linguistic diversity of languages. Specifically, (Keung et al., 2020) discovers that model’s English dev accuracy does not correlate with its performance of other languages, and (Lauscher et al., 2020) demonstrates the limitation of using English to transfer to low-resource languages. Furthermore, (Turc et al., 2021) find that other high-resource languages such as German and Russian often transfer more effectively when the set of target languages is diverse or unknown a priori. Our work builds on these findings and provides a more complete view of language relations and ZSCL performances of mLMs.

Linguistic Diversity By probing the learned representation of mLMs, (Pires et al., 2019; Limisiewicz et al., 2020) have found syntactical information implicitly encoded in layers of the multi-lingual models (typically at middle-level layers of

the architectures). (Xu et al., 2022a) shows that the pre-training and fine-tuning processes transform these features and directly impact model multi-lingual performances. Further investigation by (Lin et al., 2019) demonstrates that distances based on linguistic features, including phylogenetic and typological properties, between two languages, are correlated with cross-lingual transfer capacity. These features can be used to further improve transfer guide parameters sharing among languages (Ammar et al., 2016), or data selection (Ponti et al., 2018). Several works (Lin et al., 2019; Srinivasan et al., 2022) specifically aim to predict cross-lingual transfer performance directly, without training, only from linguistic distances of languages. Our work follows their line of reasoning but aims to address their limitation to restricted experiment settings (one-to-one transfer, model architectures, tasks, etc.). In particular, we focus on building a comprehensive picture of linguistic relations and transfer performances, in the general zero-shot mult-transfer (many-to-many) setting, for practical information extraction tasks.

Adversarial Language Learning Inspired by domain adversarial neural network (DANN) (Ganin et al., 2016) from domain adaptation research, Adversarial Language Adaptation (ALA) network can be used to extract language-invariant features useful for downstream tasks across languages. Several works have successfully adopted ALA for cross-lingual transfer setting for different tasks such as sentiment analysis (Chen et al., 2018), information extraction (Nguyen et al., 2021; Ngo Trung et al., 2021), and name tagging (Huang et al., 2019). We generalize these works to cross-lingual transfer with multiple source-target languages. This is achieved through graph-relational adversarial learning framework following (Xu et al., 2022b), a generalization of DANN.

3 Datasets

Information extraction tasks extract structured contextual information from unstructured text, thus requiring models comprehension of both syntactic and semantic knowledge in multilingual documents. In this paper, to demonstrate the heterogeneity of ZSCL for multi-lingual IE problems in practice, we experiment on two recent datasets that provide training and evaluation data in a wide range of languages.

MINION: Multi-lingual Event Detection

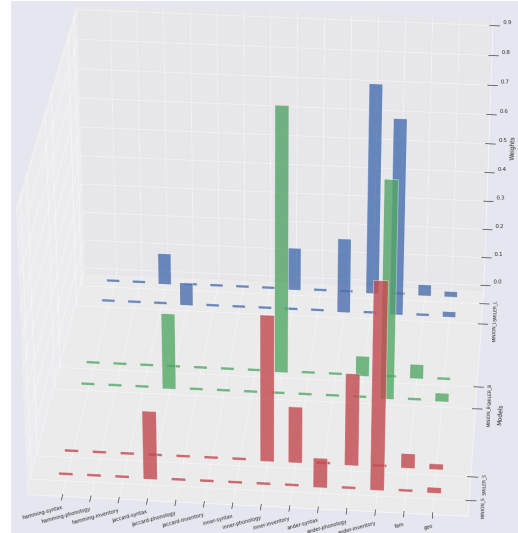


Figure 2: Feature importance weights of the optimal combined metric for each (task, model scale) setting. Small, base, and large models are represented by the colors red, green, and blue, respectively.

(MINION) (Pouan Ben Veyseh et al., 2022) annotates event triggers for 8 typologically different languages. The goal event detection task is to identify the word(s) that describe the occurrence of an event the best from a given text, also referred to as the event trigger, and classify that event into one of the 16 predetermined event types.

SMiLER: Samsung MultiLingual Entity and Relation Extraction (SMiLER) (Seganti et al., 2021) consists of annotated entities and relations from 14 languages. Given an input text, SMiLER not only requires models to identify two entity mentions in the text but also predict their relation from a set of 36 predefined relations.

The distribution of every language training set in each dataset is presented in Table 1. After categorizing the languages into high/medium/low resources groups, we notice a considerable discrepancy in the categories of shared languages between the two tasks. This reflects the diversity of actual multi-lingual data annotation processes for practical tasks. Thus, instead of balancing data across languages as in prior studies (Malkin et al., 2022), we decide to utilize the original splits of each dataset to demonstrate a realistic picture of cross-lingual transfer performance.

4 Linguistic Relations

To illustrate a comprehensive picture of the linguistic relations among the available languages, we consider different base linguistic features and how to compare them, into a total of 14 distance metrics

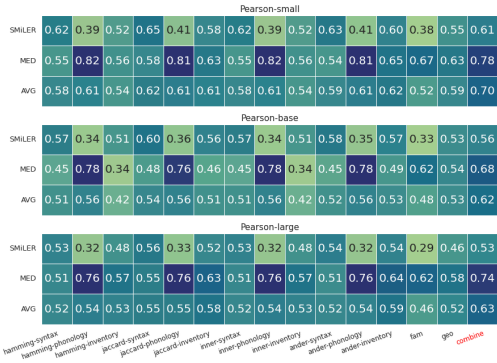


Figure 3: The language-based average Pearson correlation scores of all computed linguistic distances (including the combined metric).

4.1 Linguistic Features

Following the standard approach, We use five different linguistic features provided by the URIEL Typological Database (Littell et al., 2017), including a phylogeny feature, a geography feature, and three typological features (syntax, phonology, and inventory):

Phylogeny (fam): the membership in language families derived from the world language family tree in Glottolog (Hammarström et al., 2022)

Geography (geo): the language location based on Glottolog, more specifically the orthodromic distance between the language and a fixed point on the surface of the earth.

Syntax: the language syntactic structures derive from either WALS (Dryer and Haspelmath, 2013) or Ethnologue (Lewis, 2009)

Phonology: the phonology features extracted in a similar manner from WALS and Ethnologue

Inventory: the phonetic features derived from PHOIBLEs phonetic inventories (Moran et al., 2014)

Each of the above linguistic features is represented by a multi-dimensional binary vector for every language, where a value 0 (1) in each dimension represents the absence (presence) of a particular linguistic phenomenon for that language.

4.2 Distance Metrics

To calculate the linguistic distance between languages based on the above feature vectors, previous works only consider cosine or Euclidean distance between the binary vectors. However, numerous binary similarity measures have been proposed and play a critical role in many problems in various fields. These binary metrics are distinguished by their unique synthetic properties (negative matches, count differences, correlation, etc.),

and applying an appropriate one is the key to more accurate data analysis results. Based on the categorization in (Choi et al., 2009) which survey over 76 binary similarity measures, we decide to focus on the following 4 representative distances: Hamming, Jaccard, Inner-product, and Anderberg.

Figure 1 show the correlations between every pair of our considered linguistic distances (the full detailed distance values are provided in figure 5) in appendix A. Aside from fam and geo which are computed using Euclidean distance, each of the three typological features is computed using the chosen 4 binary metrics, resulting in a total of 14 linguistic distance metrics in the figure. We can observe significant variations in metrics based on different types of features from the correlation heatmap. In addition, there are also noticeable distinctions between several metrics within the same feature types, particularly between Anderberg and Hamming-based distances. Noted that in the context of realistic cross-lingual transfer, these two distances maybe provide more insight as transfer performance is asymmetric (Hamming) and non-zero self-distance (Anderberg). Overall, our choice of linguistic distances ensures a high diversity of correlation that can be computed from the language features.

5 Zero-shot Cross-lingual Single-transfer

To answer the first research question, we evaluate ZSCL-S scores for every language pair of each task, in three model scales: small (MiniLM (Wang et al., 2020)), base (XLM-Roberta-base (Conneau et al., 2020)), and large (XLM-Roberta-large). Next, Pearson correlations are computed between the transfer performances and linguistic distances to identify the degree to which the relations in the linguistic landscape determine a models cross-lingual transfer ability.

Experimental Setup In ZSCL-S, given a pair of source and target languages, the model is trained using labeled data from source language. The ZSCL-S score is then defined as the zero-shot evaluation of the trained model on the test set of target language.

Transfer Performance Detailed transfer scores are provided in figures 6 and 7 in appendix A. While the language-wise order of the transfer scores is maintained across different model sizes, it is not clear, however, if language identities alone are able to determine model cross-lingual transfer

ability. This is due to the significant difference between the results of the two tasks. Even more unexpectedly, model transfer scores do not increase linearly with its number of parameters

5.1 Linguistic Correlation

We determined if any of the linguistic distances defined in section 4 can explain the heterogeneity of resulting transfer performances, across all settings.

Distance-Transfer Correlation We compute the Pearson correlation between the transfer score and distance vector between each language pair. The detailed results are presented figures 8 and 9 in appendix A. While there are several distances that achieve a correlation score of over 0.7, effectively predicting the corresponding transfer performances, none of the linguistic relations are highly correlated with the transfer scores for both tasks. In particular, syntax and inventory features have above-average correlation scores for SMiLER, whereas only phonology-based distances are effective for the event detection task.

Combined Metric In order to achieve our objective of creating a universal metric that can be applied across different practical settings, we define a combined metric as a weighted average of all relevant linguistic features. For each task, the optimal weights are the solution of a simple constrained correlation linear maximization (the weights are constrained to be non-negative and sum to 1). Figure 2 compares the resulting weight importances between the two tasks across model scales. Similar to the above assessment, there is a divergence between MED and SMiLER on how the linguistic features are weighted in the optimal combined metric.

From these observations, we propose a joint combined metric that involves all three of the typological features as follows: $d_{comb} = 0.4 * d_{ander-syntax} + 0.2 * d_{inner-phonology} + 0.4 * d_{ander-inventory}$. To demonstrate the adaptability of the new distance, we provide the mean correlation scores (across all languages) of all computed linguistic distances in figure 3. Not only d_{comb} achieves the highest correlation with transfer performances overall (above 0.6 for every setting), the combined metric also greatly lessens the score’s variability amid tasks and scales of models. This implies that d_{comb} has the potential to be a general metric to approximate ZSCL perfor-

mances prior to model training. The following sections will use this combined distance for guiding the language selection and adversarial training in multi-transfer setting.

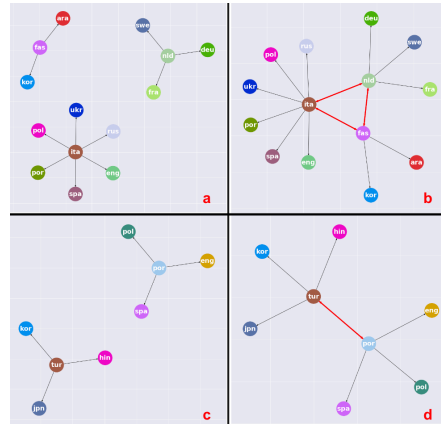


Figure 4: Language clustering results for languages in SMiLER (a and b) and MINION (c and d). The graphs on the right (b and d) are the same as the ones on the left, but with connected medoids indicated by the new red edges.

	SMALL	BASE	LARGE	MODEL_AVG	
M	medoids*	1.8	1.7	0.7	1.4
	tur*	2.7	1.0	1.0	1.5
	por*	6.0	6.1	5.8	6.0
	task_avg	3.5	2.9	2.5	3.0
S	medoids*	2.1	1.4	1.9	1.8
	ita*	3.9	3.7	3.1	3.5
	nld*	9.2	8.6	6.7	8.2
	fas*	1.8	1.7	0.7	1.4
	task_avg	4.2	3.8	3.1	3.7

Table 2: Differences in ZSCL-M scores (F1) of Inter-cluster ($medoids^*$) and Intra-cluster ($medoid_{lang}^*$) configurations over Random configuration, for tasks MED (M) and SMiLER (S).

6 Zero-shot Cross-lingual Multi-transfer

We address question Q2 by evaluating multi-transfer performances between two sets of languages. In particular, we define a transfer configuration as an experimental setting that specifies languages inside the source and target sets, and a transfer run as an actual experimental evaluation of a transfer configuration. As the number of configurations is exponential in terms of the number of languages, it is computationally impossible to evaluate every configuration, even more so for different tasks and model scales. Therefore, we focus solely on the resource-constrained scenario, which is also equivalent to the setting with the minimum number of source languages. Based on the result from the previous section, we propose to limit the configuration scope using the general combined distance metric as follows.

6.1 Language Selection

In the resource-constrained setting, our goal is to select the minimal set of source languages D_s that

can maximally transfer to a given target language set D_t . We further restrict our attention to transfer configurations with D_t as set of closely related languages in terms of cross-lingual transfer. Assuming pair-wise transfer is highly correlated with multi-transfer, these configurations can be identified by clustering languages based on the combined linguistic distance d_{comb} .

Language Clustering We use k-means (Lloyd, 1982) to partition the available languages into clusters based on d_{comb} . In particular, k-medoids¹, a variant of k-means is used for both clustering and finding the medoids (the actual data point/language that is the center of each cluster). These medoid languages should be the optimal language that transfers best to every member of its cluster. The resulting clusterings are shown in figure 4 for both MINION (4c) and SMiLER (4a). The minimal number of source languages ($|D_s|$) is chosen to be equal to the minimal number of clusters such that each cluster has at least 3 members (so we can meaningfully evaluate multi-transfer setting).

6.2 Experimental Setup

In ZSCL-M, source training set D_s consists of N_s labeled datasets, and the goal is to transfer to target cluster D_t . For MINION task, $N_s = 2$ and D_t can be the turkish cluster *tur**, or the portuguese cluster *por**. For SMiLER, $N_s = 3$ and D_t can be the italian cluster *ita**, the dutch cluster *nld**, or the farsi cluster *fas**. We are interested in verifying the following two hypotheses:

(1) Inter-Cluster Transfer: The best set of languages that transfer best across every cluster is the set of every medoid language.

(2) Intra-Cluster Transfer: The best set of languages that transfer best for a given cluster is a subset of that cluster.

6.3 Transfer Performance

Detailed results of multi-transfer performances are provided in figures 10 and 11 in appendix A, from which we can observe a clear improvement over single-transfer setting owning the additional training data. Our main interest here is how effective d_{comb} is in guiding the language selection for ZSCL-M. Table. 2 shows the differences in transfer scores of Inter-cluster (*medoids**) and Intra-cluster (*[medoid_lang]**) configura-

tions over Random configuration. Specifically, *medoids** measures inter-cluster transfer capacity of the set D_s consisting of every medoid from each cluster, to the target set D_t of all considered languages. On the other hand, *[medoid_lang]** measures intra-cluster performance of a randomly sampled set D_s of N_s languages from the corresponding cluster, to the target set D_t of every language in that cluster. Finally, Random configurations are sampled from the set of configurations that are not part of the above two configurations. Except for Inter-cluster configuration which only has one option, the results of Intra-cluster and Random configurations are the average of their sampled transfer runs.

The results from table 2 show that language selection based on d_{comb} provides a considerable boost in multi-transfer scores for every (task, configuration, model scale) setting. This suggests that these medoid languages have the potential to achieve optimal cost-performance trade-offs in multi-transfer setting. Notably between the two tasks, SMiLER has higher performance increases with only one additional source languages, despite having almost double the number of languages in target set. This implies that as the number of languages grows considerably larger, there may exist a magnitude smaller set of optimal universal languages (medoids) that are able to transfer to every language extremely well.

	SMALL		BASE		LARGE		MODEL_AVG		
	ZSCL-R	DANN	ZSCL-R	DANN	ZSCL-R	DANN	ZSCL-R	DANN	
M	medoid*	0.7	-3.7	1.4	-1.9	1.3	-2.8	1.1	-2.8
	tur*	4.1	-0.5	2.7	-2.7	3.6	-0.1	3.5	-1.1
	por*	0.6	-5.5	-0.3	-4.2	-0.3	-5.8	0.0	-5.2
	task_avg	1.8	-3.2	1.3	-2.9	1.5	-2.9	1.5	-3.0
S	medoid*	1.8	-11.8	3.2	-4.2	0.6	-1.7	1.9	-5.9
	ita*	3.0	-15.1	3.9	-5.5	2.3	-5.8	3.1	-8.8
	nld*	2.4	-10.9	2.8	-4.9	0.0	-5.1	1.7	-7.0
	fas*	-0.2	-10.6	2.1	-4.2	2.4	-2.7	1.4	-5.8
	task_avg	1.8	-12.1	3.0	-4.7	1.3	-3.8	2.0	-6.9

Table 3: Difference between transfer performances of adversarial training methods and ZSCL-M runs in inter-cluster setting.

7 Zero-shot Cross-lingual Relational-transfer

Due to limited access to multi-lingual annotators, gathering labeled data across languages is difficult. Previous section address this by careful language/data selection to optimize cost-effect. In contrast, unlabeled data is easy to collect, but leveraging it correctly for ZSCL is non-trivial. This section investigates the effectiveness of adversarial training approach for the more general ZSCL-M setting, and the possibility of further im-

¹<https://en.wikipedia.org/wiki/K-medoids>

proving multi-lingual transfer through explicitly integrating our transfer-correlated linguistic relations.

7.1 Experimental Setup

We follow the same setup as in ZSCL-M, but each language is accompanied by an unlabeled dataset which can be used for training. The model use labeled data from the source cluster to learn the task, whereas unlabeled data from another cluster is used to help transfer source performance to that target cluster. The goal is to bridge the performance between 2 different language clusters with the aid of given unlabeled text.

Adversarial Language Adaptation A typical method use for ZSCL-S is adversarial language adaptation (ALA) (Chen et al., 2018; Nguyen et al., 2021) which employs a language discriminator that takes an encoded representation from mLMs as its input and predicts its origin (language). By pushing the encoder to both minimize the downstream loss and maximally misdirect the language predictor (adversarial training), the resulting representation can be indiscriminate with respect to the shift between the languages while also discriminative for the main learning task. Apply ALA for ZSCL-M setting is equivalent to applying DANN for a single joint source domain to a single joint target domain (the union of every languages in D_s and D_t , respectively).

Zero-shot Cross-lingual Relational-transfer We extend ALA to the case of multiple source and target languages through Graph-relational domain adaptation (GrDA) (Xu et al., 2022b), a generalization of DANN to multi-domains adaptation setting by introducing a domain graph that captures heterogeneous relations among domains. GrDA relaxes the strict uniform alignment of DANN to allow flexible and effective adaptation between distant domains. We use the language clustering graphs on the right of figure 4 as domain graphs for GrDA. Noted that additional direct connections between medoid languages are introduced (red edges) to facilitate inter-cluster transfer. We refer to this adversarial learning process for ZSCL that directly embeds the linguistic relations into the representations as Zero-shot Cross-lingual Relational-transfer (ZSCL-R).

7.2 Transfer Performance

Performances of the baseline DANN and the proposed adversarial training method ZSCL-R are

provided in table. 3, which follows the same format as table. 2. However, instead of comparing against Random configurations, they are compared directly with results of ZSCL-M runs of the corresponding configurations, but only in terms of inter-cluster transfer. The negative results of DANN confirm that strictly aligning language representations uniformly is not effective in ZSCL-M setting. As the model scale gets smaller, model’s representation becomes less expressive, whereas the language-invariant feature of all languages is harder to capture as the number of languages grows. Thus, the adverse effect gets significantly worse on small models for SMiLER task (-12 points on average). In contrast, ZSCL-R provides consistent improvements over ZSCL-M for most configurations. Due to the flexibility of GrDA alignment, ZSCL-R performs even better as the number of languages increases, effectively leveraging the additional unlabeled data to help improve inter-cluster transfer ability of models.

8 Conclusion

We explored the general cross-lingual transfer learning setting where multiple source and target languages are involved. Our experiments on two practical information extraction tasks across different model scales and languages reveal new general insights on cross-lingual transfer learning: (1) There is a correlation between linguistic distances and single-transfer performances; however, simplistic measures based on syntax features are only sufficient for syntactic-based tasks. We develop a combined distance based on various metrics and linguistic features that achieves a high correlation with cross-lingual transfer score robustly across all settings. (2) The proposed combined metric provides useful directions for language clustering and selection to achieve optimal cost-performance trade-off in multi-transfer to a specific group of languages. (3) Finally, linguistic relations can be leveraged with unlabeled data for adversarial training to help generalize to a new group of languages with minimum additional annotation cost. Our findings collectively suggest multi-transfer as a new baseline for cross-lingual learning, and provide a baseline for efficient and effective multi-transfer together with promising directions that future work can further improve upon.

Limitations

Compared to prior cross-lingual transfer papers (Srinivasan et al., 2022), our work aims to demonstrate generalization across various hyperparameters and design choices that affect the results of previous investigations on the topic. Consequently, this has led to significant computational demands, forcing us to limit and simplify some aspects of our experiments to ensure manageability. This section outlines what we have and has not been able to address, and suggests promising future directions that can be followed from our findings.

First, our combined metric is heuristically defined based on the transfer-distance correlation scores. Although the optimal metric for all situations is impossible to find and likely non-existent, we anticipate that the ideal metric won't differ substantially across settings of tasks and model architectures, and may only vary slightly from ours. Nevertheless, an in-depth analysis is needed on how a significant change in the distance metric can impact ZSML-M and particularly ZSML-R, which is explicitly guided by the metric.

Second, We only experiment with a minimum number of language clusters needed for an effective evaluation of multi-transfer. There is no guarantee that this minimal cost strategy is also the one with the best cost-performance trade-off, which can depend on the number of languages and task availability. Future work may investigate this trade-off as the problem scales to hundreds of languages, in particular validating the hypothesis stated in section 6.3: the best trade-off point (the number of clusters) is an order of magnitude smaller than the number of available languages.

Third, the decision to connect the medoids in the language cluster graph for ZSCL-R is the simplest solution to create a connected language-relation graph. This, however, also implies the path between any two languages has a maximum length of 3, which does not reflect the actual relation distance among languages. Further research is needed to ascertain the optimal language relational-graph, especially in intricate scenarios that involve many languages and tasks.

Finally, our model scale only stops at hundreds of millions of parameters, which are no longer considered large scale by today's standard. Further experimentation can test if our results hold for the current billion scale models. A more interesting direction would be to investigate cross-lingual multi-

transfer performance of parameter-efficient tuning (Chen et al., 2022) and instruction tuning (Wei et al., 2022), which have become the standard approaches for fine-tuning these massive scale models.

References

- Waleed Ammar, George Mulcaire, Miguel Ballesteros, Chris Dyer, and Noah A. Smith. 2016. [Many languages, one parser](#). *Transactions of the Association for Computational Linguistics*, 4:431–444.
- Guanzheng Chen, Fangyu Liu, Zaiqiao Meng, and Shangsong Liang. 2022. [Revisiting parameter-efficient tuning: Are we really there yet?](#) In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2612–2626, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Xilun Chen, Yu Sun, Ben Athiwaratkun, Claire Cardie, and Kilian Weinberger. 2018. [Adversarial deep averaging networks for cross-lingual sentiment classification](#). *Transactions of the Association for Computational Linguistics*, 6:557–570.
- Zewen Chi, Li Dong, Furu Wei, Nan Yang, Saksham Singhal, Wenhui Wang, Xia Song, Xian-Ling Mao, Heyan Huang, and Ming Zhou. 2021. [InfoXLM: An information-theoretic framework for cross-lingual language model pre-training](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics.
- SHC Choi, Sung-Hyuk Cha, and Charles Tappert. 2009. A survey of binary similarity and distance measures. *J. Syst. Cybern. Inf.*, 8.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Baej Dolicki and Gerasimos Spanakis. 2021. [Analysing the impact of linguistic features on cross-lingual transfer](#).
- Matthew S. Dryer and Martin Haspelmath, editors. 2013. *WALS Online (v2020.3)*. Zenodo.
- Yuwei Fang, Shuohang Wang, Zhe Gan, Siqi Sun, and Jingjing Liu. 2021. [FILTER: An enhanced fusion method for cross-lingual language understanding](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(14):12776–12784.

- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario March, and Victor Lempitsky. 2016. [Domain-adversarial training of neural networks](#). *Journal of Machine Learning Research*, 17(59):1–35.
- Naman Goyal, Jingfei Du, Myle Ott, Giri Anantharaman, and Alexis Conneau. 2021. [Larger-scale transformers for multilingual masked language modeling](#). In *Proceedings of the 6th Workshop on Representation Learning for NLP (RepL4NLP-2021)*. Association for Computational Linguistics.
- Harald Hammarström, Robert Forkel, Martin Haspelmath, and Sebastian Bank. 2022. [glottolog/glottolog: Glottolog database 4.7](#).
- Lifu Huang, Heng Ji, and Jonathan May. 2019. [Cross-lingual multi-level adversarial transfer to enhance low-resource name tagging](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3823–3833, Minneapolis, Minnesota. Association for Computational Linguistics.
- Phillip Keung, Yichao Lu, Julian Salazar, and Vikas Bhardwaj. 2020. [Don’t use English dev: On the zero-shot cross-lingual evaluation of contextual embeddings](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 549–554, Online. Association for Computational Linguistics.
- Anne Lauscher, Vinit Ravishankar, Ivan Vulić, and Goran Glavaš. 2020. [From zero to hero: On the limitations of zero-shot language transfer with multilingual Transformers](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4483–4499, Online. Association for Computational Linguistics.
- M. Paul Lewis, editor. 2009. *Ethnologue: Languages of the World*, sixteenth edition. SIL International, Dallas, TX, USA.
- Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun Shou, Daxin Jiang, Guihong Cao, Xiaodong Fan, Ruofei Zhang, Rahul Agrawal, Edward Cui, Sining Wei, Taroon Bharti, Ying Qiao, Jiun-Hung Chen, Winnie Wu, Shuguang Liu, Fan Yang, Daniel Campos, Rangan Majumder, and Ming Zhou. 2020. [XGLUE: A new benchmark dataset for cross-lingual pre-training, understanding and generation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.
- Tomasz Limisiewicz, David Mareček, and Rudolf Rosa. 2020. [Universal Dependencies According to BERT: Both More Specific and More General](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2710–2722, Online. Association for Computational Linguistics.
- Yu-Hsiang Lin, Chian-Yu Chen, Jean Lee, Zirui Li, Yuyan Zhang, Mengzhou Xia, Shruti Rijhwani, Junxian He, Zhisong Zhang, Xuezhe Ma, Antonios Anastasopoulos, Patrick Littell, and Graham Neubig. 2019. [Choosing transfer languages for cross-lingual learning](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3125–3135, Florence, Italy. Association for Computational Linguistics.
- Patrick Littell, David R. Mortensen, Ke Lin, Katherine Kairis, Carlisle Turner, and Lori Levin. 2017. [URIEL and lang2vec: Representing languages as typological, geographical, and phylogenetic vectors](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 8–14, Valencia, Spain. Association for Computational Linguistics.
- S. Lloyd. 1982. [Least squares quantization in pcm](#). *IEEE Transactions on Information Theory*, 28(2):129–137.
- Dan Malkin, Tomasz Limisiewicz, and Gabriel Stanovsky. 2022. [A balanced data approach for evaluating cross-lingual transfer: Mapping the linguistic blood bank](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics.
- Steven Moran, Daniel McCloy, and Richard Wright, editors. 2014. *PHOIBLE Online*. Max Planck Institute for Evolutionary Anthropology, Leipzig.
- Nghia Ngo Trung, Duy Phung, and Thien Huu Nguyen. 2021. [Unsupervised domain adaptation for event detection using domain-specific adapters](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4015–4025, Online. Association for Computational Linguistics.
- Minh Van Nguyen, Tuan Ngo Nguyen, Bonan Min, and Thien Huu Nguyen. 2021. [Crosslingual transfer learning for relation and event extraction via word category and class alignments](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5414–5426, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jason Phang, Iacer Calixto, Phu Mon Htut, Yada Pruksachatkun, Haokun Liu, Clara Vania, Katharina Kann, and Samuel R. Bowman. 2020. [English intermediate-task training improves zero-shot cross-lingual transfer too](#). In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 557–575, Suzhou, China. Association for Computational Linguistics.

- Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. [How multilingual is multilingual BERT?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4996–5001, Florence, Italy. Association for Computational Linguistics.
- Edoardo Maria Ponti, Roi Reichart, Anna Korhonen, and Ivan Vulić. 2018. [Isomorphic transfer of syntactic structures in cross-lingual NLP.](#) In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1531–1542, Melbourne, Australia. Association for Computational Linguistics.
- Amir Pouran Ben Veyseh, Minh Van Nguyen, Franck Dernoncourt, and Thien Nguyen. 2022. [MINION: a large-scale and diverse dataset for multilingual event detection.](#) In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2286–2299, Seattle, United States. Association for Computational Linguistics.
- Sebastian Ruder, Noah Constant, Jan Botha, Aditya Siddhant, Orhan Firat, Jinlan Fu, Pengfei Liu, Junjie Hu, Dan Garrette, Graham Neubig, and Melvin Johnson. 2021. [Xtreme-r: Towards more challenging and nuanced multilingual evaluation.](#) *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.
- Alessandro Seganti, Klaudia Firląg, Helena Skowronska, Michał Satława, and Piotr Andruszkiewicz. 2021. [Multilingual entity and relation extraction dataset and model.](#) In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1946–1955, Online. Association for Computational Linguistics.
- Anirudh Srinivasan, Gauri Kholkar, Rahul Kejriwal, Tanuja Ganu, Sandipan Dandapat, Sunayana Sitaram, Balakrishnan Santhanam, Somak Aditya, Kalika Bali, and Monojit Choudhury. 2022. [Litmus predictor: An ai assistant for building reliable, high-performing and fair multilingual nlp systems.](#) In *Thirty-sixth AAAI Conference on Artificial Intelligence*. AAAI. System Demonstration.
- Anirudh Srinivasan, Sunayana Sitaram, Tanuja Ganu, Sandipan Dandapat, Kalika Bali, and Monojit Choudhury. 2021. Predicting the performance of multilingual nlp models. *ArXiv*, abs/2110.08875.
- Iulia Turc, Kenton Lee, Jacob Eisenstein, Ming-Wei Chang, and Kristina Toutanova. 2021. [Revisiting the primacy of english in zero-shot cross-lingual transfer.](#)
- Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS’20*, Red Hook, NY, USA. Curran Associates Inc.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022. [Finetuned language models are zero-shot learners.](#) In *International Conference on Learning Representations*.
- Ningyu Xu, Tao Gui, Ruotian Ma, Qi Zhang, Jingting Ye, Menghan Zhang, and Xuanjing Huang. 2022a. [Cross-linguistic syntactic difference in multilingual BERT: How good is it and how does it affect transfer?](#) In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8073–8092, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Zihao Xu, Hao He, Guang-He Lee, Bernie Wang, and Hao Wang. 2022b. [Graph-relational domain adaptation.](#) In *International Conference on Learning Representations*.

A Detailed Experimental Results

In this section, we provide detailed result values for our experiments, including: linguistic distances in figure 5, ZSCL-S scores in figures 6 and 7, transfer-distance correlations in figures 8 and 9, ZSCL-M scores in figures 10 and 11, DANN scores in figures 12 and 13, and finally ZSCL-R scores in figures 14 and 15.

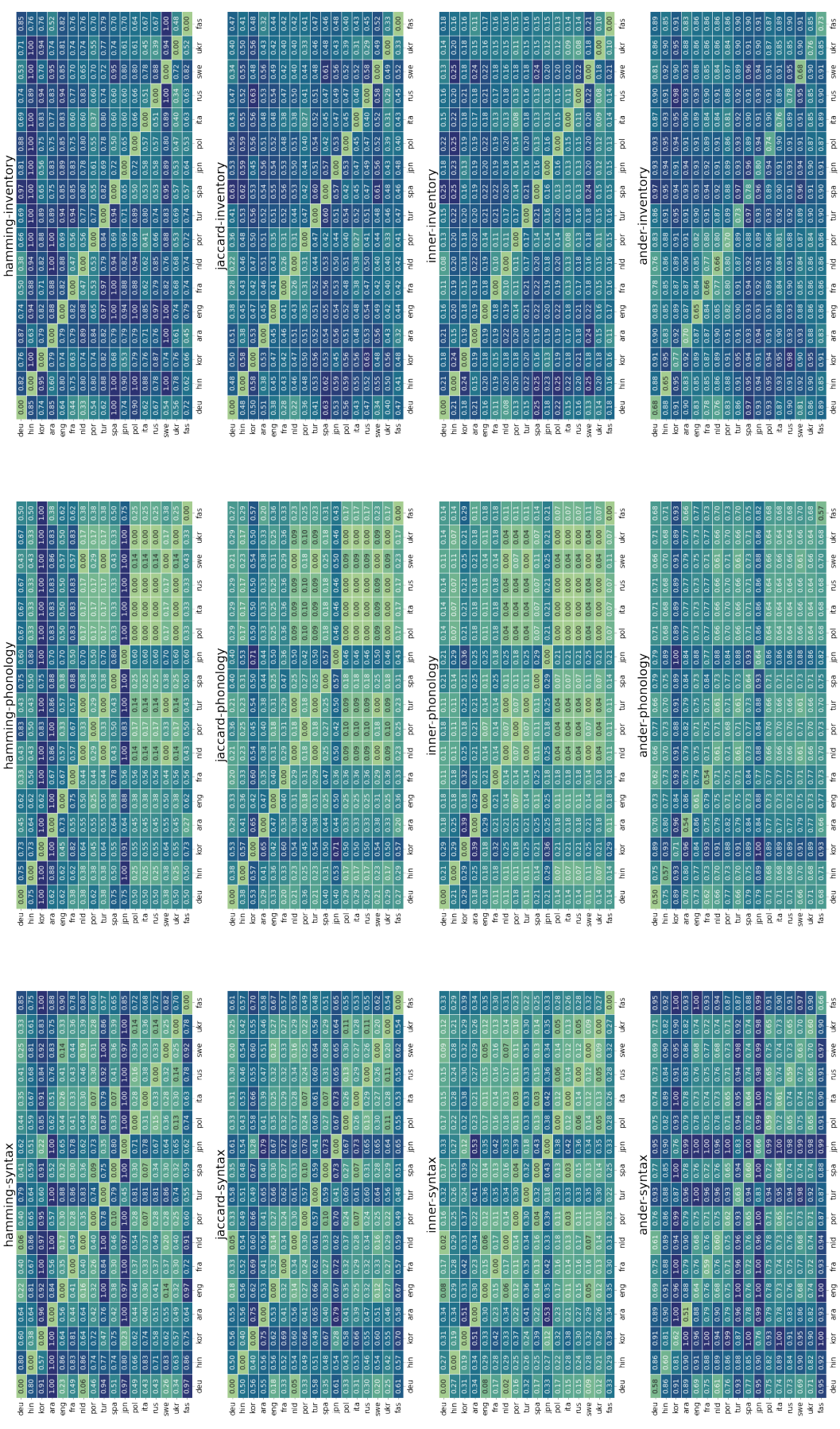


Figure 5: Detailed linguistic distances



Figure 6: Detailed transfer performances for MINION task in ZSCL-S setting.

ZSCL-S - SMiLER - small

eng	84.60	37.70	53.80	51.80	30.40	55.70	54.90	51.10	58.10	48.60	37.30	47.70	33.10	58.60	50.24
ara	50.80	95.60	57.70	54.70	59.20	65.20	67.20	49.30	63.90	67.50	69.60	55.40	63.50	65.90	63.25
deu	66.20	42.60	90.00	77.50	36.60	80.20	78.60	63.70	81.40	74.00	41.50	64.60	36.00	78.50	65.10
spa	64.90	39.30	74.80	84.60	35.30	79.10	80.20	56.40	78.20	71.40	38.70	58.20	35.00	82.40	62.75
fas	51.40	71.20	61.30	62.00	87.80	64.80	54.20	42.70	55.70	60.60	61.80	47.20	51.60	64.20	59.75
fra	64.80	32.90	76.20	75.80	27.20	88.10	78.00	55.90	74.80	67.60	34.90	59.90	29.50	78.90	60.32
ita	70.80	35.60	78.40	79.40	31.00	85.30	95.50	57.70	80.10	72.20	36.40	62.10	30.90	82.00	64.10
kor	67.10	54.10	73.40	68.10	51.10	72.80	69.90	91.80	75.30	65.70	51.60	58.80	47.70	72.10	65.68
nld	70.50	55.40	85.90	75.20	44.00	85.70	84.80	53.90	93.80	81.30	55.40	69.40	50.30	87.30	70.92
pol	66.20	53.40	79.50	77.00	42.80	80.10	78.70	55.90	80.50	93.00	55.60	68.40	50.30	80.50	68.71
rus	75.30	81.70	82.70	77.80	65.00	83.40	81.40	60.60	85.20	83.60	94.20	74.80	85.10	84.80	79.69
swe	56.30	44.70	74.50	65.80	32.30	72.60	72.50	48.70	91.00	70.50	40.30	93.90	28.80	78.00	62.14
ukr	60.50	64.00	70.00	68.70	55.10	69.20	71.30	57.50	70.50	84.70	82.80	62.60	85.60	72.70	69.66
por	65.50	37.20	75.10	79.70	30.60	79.20	80.40	57.60	78.80	70.70	37.90	59.30	31.80	89.80	62.40
avg	65.35	53.24	73.81	71.29	44.89	75.81	74.83	57.34	76.24	72.24	52.71	63.02	47.09	76.84	64.62
	eng	ara	deu	spa	fas	fra	ita	kor	nld	pol	rus	swe	ukr	por	avg

ZSCL-S - SMiLER - base

eng	84.70	42.20	56.70	54.00	33.80	57.80	56.20	54.10	59.10	50.10	38.50	48.40	32.20	59.80	51.97
ara	53.80	95.00	57.70	61.90	68.10	66.70	66.50	54.70	63.40	63.90	71.00	56.40	67.00	64.60	65.05
deu	66.60	43.40	90.10	79.50	36.70	81.30	79.60	66.50	82.40	74.60	43.20	66.50	36.90	79.70	66.21
spa	65.70	41.90	76.50	84.00	36.40	80.00	78.70	60.40	79.10	70.00	38.50	59.30	32.40	82.40	63.24
fas	50.40	72.40	65.10	66.20	88.60	64.40	55.10	52.70	59.40	69.50	68.70	52.20	55.10	63.90	63.12
fra	62.00	33.50	76.90	76.00	28.10	88.30	76.70	56.70	76.00	68.00	34.70	58.30	28.80	79.60	60.26
ita	69.30	36.50	79.50	80.60	31.60	85.70	92.10	60.40	80.10	71.90	37.00	65.60	29.90	82.90	64.51
kor	69.40	56.10	75.60	69.90	51.20	75.40	68.10	92.60	76.90	66.30	52.80	61.90	48.30	71.30	66.84
nld	74.10	59.00	87.30	80.90	46.90	86.10	84.30	66.60	94.50	81.70	56.00	73.50	46.40	86.80	73.15
pol	66.50	55.00	79.70	79.40	46.80	80.60	78.00	63.10	80.30	93.60	56.70	70.00	50.40	82.50	70.19
rus	72.10	82.40	81.70	82.50	69.20	87.50	81.70	77.10	85.80	84.00	97.50	78.70	84.70	85.10	82.14
swe	58.50	49.30	76.00	67.20	38.90	77.20	73.10	49.70	84.90	71.10	41.70	94.70	31.50	80.60	63.89
ukr	58.70	68.70	71.50	74.90	63.90	75.60	73.60	60.10	78.00	84.70	87.50	73.60	80.20	77.70	73.48
por	65.70	40.70	77.20	80.40	32.40	81.00	80.50	60.80	80.00	71.10	37.50	65.30	31.70	88.50	63.77
avg	65.54	55.44	75.11	74.10	48.04	77.69	74.59	62.54	77.14	72.89	54.38	66.03	46.82	77.53	66.27
	eng	ara	deu	spa	fas	fra	ita	kor	nld	pol	rus	swe	ukr	por	avg

ZSCL-S - SMiLER - large

eng	85.90	41.00	58.00	55.30	37.10	58.20	57.60	58.60	60.70	52.70	41.30	51.20	35.70	62.00	53.95
ara	59.00	95.90	68.30	67.70	71.70	70.60	73.40	63.30	66.60	73.70	77.70	63.90	71.40	75.40	71.33
deu	67.90	45.80	91.10	82.00	40.40	83.30	82.00	71.20	84.80	77.40	46.00	71.30	40.20	83.50	69.06
spa	67.10	42.00	78.30	86.50	41.10	82.60	83.80	64.20	80.80	74.70	40.70	65.70	35.70	86.30	66.39
fas	58.60	81.60	65.80	76.00	89.80	74.90	69.70	64.40	67.80	77.60	77.20	60.00	68.80	80.40	72.33
fra	68.30	37.90	81.70	81.20	31.50	88.90	82.80	61.60	79.80	71.10	36.90	67.60	32.10	83.30	64.62
ita	72.60	38.00	83.40	84.90	36.30	87.60	95.50	64.60	82.90	74.60	37.60	72.20	33.10	86.50	67.84
kor	77.80	59.50	80.40	76.00	53.80	78.10	78.20	92.50	80.40	75.50	56.40	69.70	53.30	79.40	72.21
nld	73.20	58.40	88.70	84.50	51.70	88.80	87.40	73.60	94.50	84.70	59.00	78.10	53.20	89.50	76.09
pol	68.80	58.10	82.20	81.30	51.50	83.00	82.40	71.10	81.50	94.60	59.90	74.90	54.10	84.20	73.40
rus	74.90	88.10	87.90	86.10	74.30	91.70	88.20	76.50	88.80	86.00	96.80	80.80	88.50	89.50	85.58
swe	61.10	55.20	77.70	69.10	44.60	78.90	74.90	56.50	91.90	76.70	48.50	95.00	34.90	86.30	67.95
ukr	59.70	81.30	79.60	81.70	72.50	85.60	76.40	69.50	83.00	89.30	89.30	70.70	88.90	84.00	79.39
por	67.90	41.60	80.80	83.90	36.80	83.50	84.20	64.30	81.90	73.70	40.60	71.10	35.50	89.80	66.83
avg	68.77	58.89	78.85	78.30	52.36	81.12	79.75	67.99	80.39	77.31	57.71	70.87	51.81	82.86	70.50
	eng	ara	deu	spa	fas	fra	ita	kor	nld	pol	rus	swe	ukr	por	avg

Figure 7: Detailed transfer performances for SMiLER task in ZSCL-S setting.

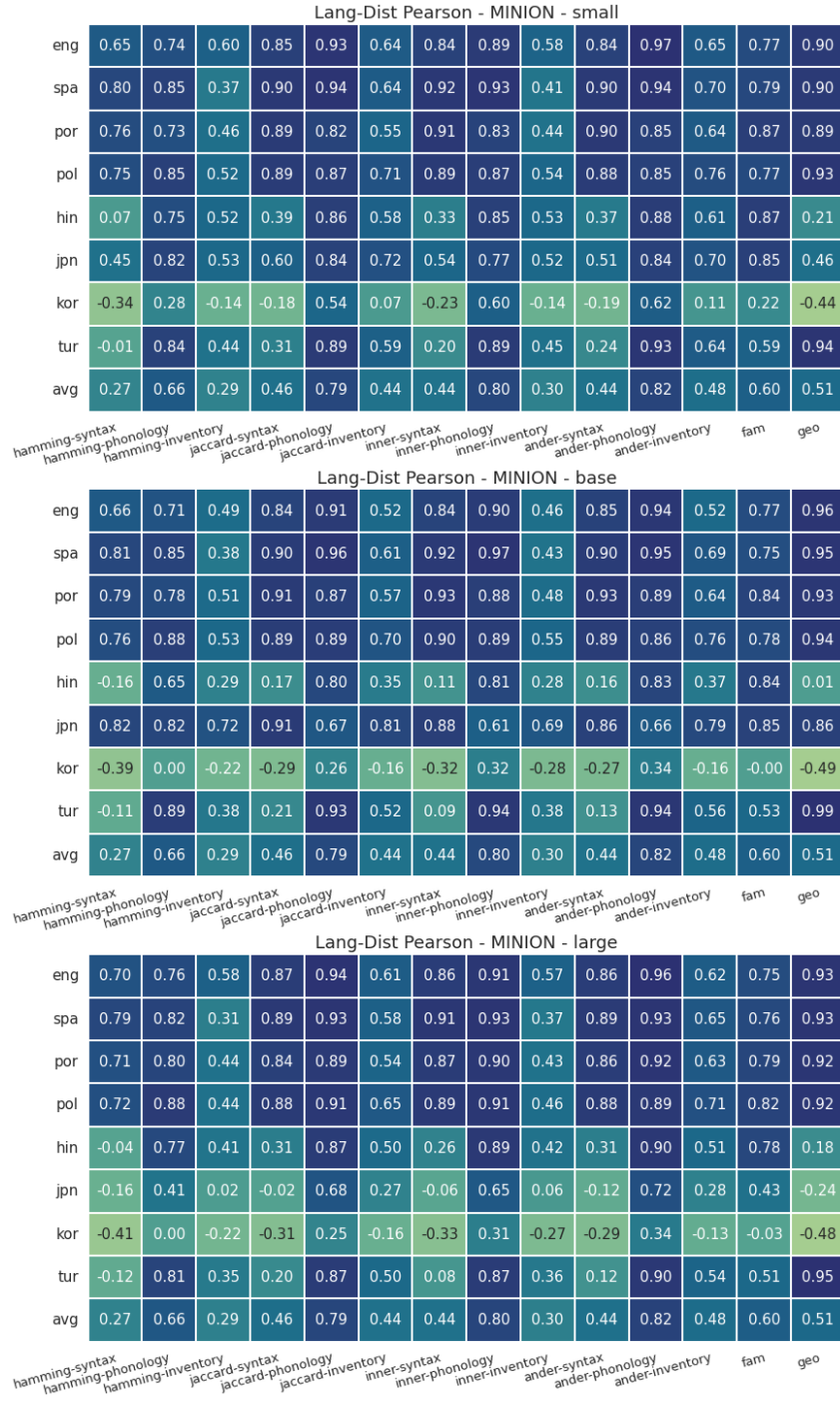


Figure 8: Detailed Pearson correlations between ZSCL-S transfer scores and linguistic distances for MINION task.

Lang-Dist Pearson - SMiLER - small

eng	0.60	0.72	0.53	0.59	0.71	0.50	0.59	0.71	0.53	0.56	0.68	0.45	0.30	0.30
ara	0.29	0.65	0.73	0.38	0.60	0.68	0.30	0.60	0.73	0.37	0.57	0.61	0.64	0.20
deu	0.55	0.28	0.59	0.55	0.41	0.45	0.55	0.48	0.50	0.57	0.33	0.37	0.34	0.26
spa	0.51	0.26	0.48	0.56	0.39	0.47	0.52	0.39	0.47	0.52	0.34	0.46	0.43	0.33
fas	0.63	0.70	0.84	0.66	0.49	0.77	0.58	0.46	0.82	0.64	0.41	0.66	-0.25	0.89
fra	0.44	0.25	0.60	0.46	0.41	0.43	0.43	0.47	0.51	0.46	0.35	0.36	0.31	0.28
ita	0.60	0.39	0.74	0.61	0.40	0.62	0.56	0.38	0.64	0.60	0.32	0.58	0.24	0.41
kor	0.91	0.79	0.77	0.92	0.87	0.75	0.85	0.76	0.79	0.90	0.85	0.60	0.88	0.70
nld	0.56	0.68	0.58	0.54	0.58	0.42	0.54	0.59	0.46	0.57	0.50	0.37	0.25	0.30
pol	0.68	0.61	0.61	0.72	0.60	0.58	0.67	0.56	0.60	0.70	0.52	0.52	0.22	0.35
rus	0.34	0.38	0.65	0.34	0.31	0.67	0.24	0.27	0.65	0.29	0.25	0.58	-0.18	0.73
swe	0.58	0.72	0.78	0.59	0.63	0.70	0.55	0.61	0.75	0.57	0.56	0.64	0.21	0.46
ukr	0.27	0.37	0.67	0.27	0.29	0.61	0.17	0.25	0.65	0.23	0.22	0.47	-0.19	-0.06
por	0.54	0.39	0.66	0.58	0.44	0.51	0.54	0.44	0.57	0.56	0.37	0.43	0.37	0.36
avg	0.40	0.46	0.60	0.42	0.45	0.52	0.38	0.45	0.58	0.41	0.39	0.44	0.19	0.28

hamming-syntax hamming-phonology hamming-inventory jaccard-syntax jaccard-phonology jaccard-inventory inner-syntax inner-phonology inner-inventory ander-syntax ander-phonology ander-inventory fam geo

Lang-Dist Pearson - SMiLER - base

eng	0.64	0.72	0.56	0.61	0.67	0.51	0.59	0.67	0.52	0.56	0.63	0.47	0.19	0.30
ara	0.26	0.61	0.69	0.36	0.58	0.66	0.28	0.57	0.71	0.34	0.54	0.59	0.62	0.19
deu	0.55	0.26	0.59	0.55	0.39	0.46	0.56	0.46	0.52	0.56	0.31	0.39	0.37	0.25
spa	0.43	0.26	0.46	0.50	0.40	0.43	0.46	0.39	0.46	0.48	0.35	0.42	0.39	0.28
fas	0.56	0.72	0.80	0.60	0.51	0.74	0.51	0.49	0.79	0.58	0.44	0.64	-0.27	0.87
fra	0.44	0.19	0.55	0.48	0.37	0.39	0.44	0.44	0.48	0.47	0.31	0.31	0.31	0.26
ita	0.60	0.43	0.73	0.63	0.45	0.62	0.59	0.43	0.64	0.62	0.38	0.59	0.30	0.45
kor	0.83	0.65	0.61	0.85	0.78	0.59	0.81	0.70	0.64	0.86	0.74	0.41	0.78	0.77
nld	0.54	0.68	0.63	0.54	0.59	0.47	0.54	0.61	0.54	0.57	0.50	0.41	0.29	0.28
pol	0.59	0.65	0.65	0.66	0.64	0.62	0.62	0.61	0.65	0.63	0.58	0.58	0.31	0.31
rus	0.30	0.37	0.64	0.31	0.31	0.66	0.21	0.27	0.64	0.26	0.26	0.57	-0.16	0.75
swe	0.55	0.79	0.70	0.58	0.70	0.64	0.54	0.69	0.72	0.55	0.59	0.57	0.25	0.43
ukr	0.20	0.32	0.63	0.20	0.25	0.56	0.10	0.21	0.60	0.16	0.19	0.42	-0.24	-0.11
por	0.57	0.43	0.65	0.62	0.50	0.50	0.59	0.50	0.56	0.60	0.44	0.42	0.45	0.40
avg	0.40	0.46	0.60	0.42	0.45	0.52	0.38	0.45	0.58	0.41	0.39	0.44	0.19	0.28

hamming-syntax hamming-phonology hamming-inventory jaccard-syntax jaccard-phonology jaccard-inventory inner-syntax inner-phonology inner-inventory ander-syntax ander-phonology ander-inventory fam geo

Lang-Dist Pearson - SMiLER - large

eng	0.43	0.69	0.57	0.38	0.54	0.51	0.36	0.56	0.53	0.32	0.48	0.46	0.02	0.07
ara	0.20	0.58	0.65	0.29	0.56	0.61	0.22	0.56	0.68	0.28	0.52	0.54	0.53	0.15
deu	0.47	0.22	0.53	0.45	0.34	0.39	0.45	0.40	0.46	0.49	0.26	0.31	0.21	0.19
spa	0.32	0.19	0.46	0.37	0.31	0.41	0.33	0.30	0.49	0.35	0.26	0.41	0.31	0.15
fas	0.52	0.72	0.78	0.57	0.53	0.71	0.49	0.51	0.77	0.56	0.45	0.61	-0.25	0.86
fra	0.32	0.11	0.45	0.37	0.33	0.31	0.35	0.40	0.43	0.37	0.27	0.22	0.33	0.13
ita	0.43	0.38	0.75	0.44	0.37	0.60	0.39	0.35	0.67	0.42	0.28	0.53	0.18	0.22
kor	0.82	0.60	0.61	0.85	0.73	0.60	0.80	0.63	0.65	0.87	0.68	0.42	0.78	0.79
nld	0.46	0.72	0.60	0.47	0.62	0.44	0.48	0.64	0.52	0.49	0.54	0.38	0.31	0.20
pol	0.47	0.57	0.64	0.51	0.53	0.57	0.45	0.50	0.62	0.47	0.46	0.51	0.12	0.14
rus	0.22	0.34	0.57	0.23	0.27	0.59	0.13	0.23	0.57	0.17	0.24	0.50	-0.20	0.74
swe	0.45	0.73	0.73	0.46	0.63	0.64	0.43	0.61	0.72	0.44	0.53	0.58	0.17	0.34
ukr	0.13	0.32	0.64	0.14	0.25	0.58	0.05	0.21	0.61	0.10	0.19	0.44	-0.22	-0.12
por	0.35	0.30	0.48	0.40	0.37	0.33	0.37	0.37	0.42	0.38	0.31	0.23	0.32	0.13
avg	0.40	0.46	0.60	0.42	0.45	0.52	0.38	0.45	0.58	0.41	0.39	0.44	0.19	0.28

hamming-syntax hamming-phonology hamming-inventory jaccard-syntax jaccard-phonology jaccard-inventory inner-syntax inner-phonology inner-inventory ander-syntax ander-phonology ander-inventory fam geo

Figure 9: Detailed Pearson correlations between ZSCL-S transfer scores and linguistic distances for SMiLER task.

ZSCL-M - MINION - small											
	english	hindi	japanese	korean	polish	portuguese	spanish	turkey	average	tur_ave	por_ave
mediator	71.45	61.06	52.26	53.82	67.65	64.55	68.56	76.48	64.48	60.91	68.05
tur	62.36	62.35	52.72	57.33	60.56	52.40	60.53	64.20	59.06	59.15	58.97
por	86.72	62.02	39.83	55.20	68.66	67.68	76.28	66.65	65.38	55.92	74.83
avg	70.91	60.85	47.14	55.11	69.08	62.85	72.39	62.88	62.65	56.50	68.81
ZSCL-M - MINION - base											
	english	hindi	japanese	korean	polish	portuguese	spanish	turkey	average	tur_ave	por_ave
mediator	74.80	65.34	57.55	59.63	76.74	70.76	72.94	78.00	69.47	65.13	73.81
tur	67.67	66.22	54.02	58.27	72.59	58.56	68.69	70.97	64.63	62.37	66.88
por	88.20	67.47	50.27	59.81	79.72	72.39	80.55	72.45	71.36	62.50	80.22
avg	74.51	66.32	52.02	58.51	76.12	68.28	77.65	68.74	67.77	61.40	74.14
ZSCL-M - MINION - large											
	english	hindi	japanese	korean	polish	portuguese	spanish	turkey	average	tur_ave	por_ave
mediator	75.28	63.30	47.52	57.41	75.36	69.09	73.03	76.05	67.13	61.07	73.19
tur	67.23	63.89	55.94	55.77	65.05	57.45	66.17	68.20	62.46	60.95	63.97
por	87.71	62.97	49.05	58.78	74.36	72.55	80.19	68.53	69.27	59.83	78.70
avg	74.88	63.86	51.38	57.72	73.99	66.77	75.90	66.75	66.41	59.93	72.88

Figure 10: Detailed transfer performances for MINION task in ZSCL-M setting.

ZSCL-M - SMILER - small																		
	en	ar	de	es	fa	fr	it	ko	nl	pl	ru	sv	uk	pt	average	ita_avg	nld_avg	fas_avg
mediator	59.96	73.19	82.82	81.66	90.44	80.69	94.37	71.97	94.54	81.89	89.40	89.95	80.34	82.98	82.44	81.51	87.00	78.53
tur	65.00	75.88	79.12	82.90	73.01	78.44	85.65	71.83	86.03	89.78	93.62	75.85	88.34	83.43	80.63	84.10	79.86	73.57
por	59.83	69.33	90.78	82.12	66.22	88.40	85.82	77.79	94.59	83.74	87.37	91.98	76.52	82.03	81.18	79.63	91.44	71.11
avg	50.43	95.33	60.56	52.28	89.82	44.31	49.59	90.16	68.22	60.40	85.26	51.55	71.79	48.94	65.62	59.81	56.16	91.77
random	57.45	77.32	82.23	77.69	76.89	77.86	80.46	80.58	87.19	88.36	92.84	81.76	85.07	79.81	80.39	80.24	82.26	78.26
ZSCL-M - SMILER - base																		
	en	ar	de	es	fa	fr	it	ko	nl	pl	ru	sv	uk	pt	average	ita_avg	nld_avg	fas_avg
mediator	61.43	78.96	84.36	81.97	91.37	80.93	93.76	77.02	94.82	83.11	88.77	89.93	79.15	84.39	83.57	81.80	87.51	82.45
tur	65.46	72.55	81.01	83.62	74.10	79.00	86.56	73.61	86.79	91.23	94.50	77.10	93.53	85.00	81.72	85.70	80.97	73.42
por	61.42	68.56	90.85	83.21	71.63	88.89	87.29	79.41	94.75	84.54	89.34	93.15	78.94	84.10	82.58	81.26	91.91	73.20
avg	56.51	95.28	62.90	57.34	91.37	45.27	50.40	91.95	71.91	64.02	85.84	49.58	80.99	53.84	68.37	64.13	57.41	92.87
random	59.74	78.09	82.95	79.68	80.98	78.93	81.68	84.05	87.72	89.00	93.23	83.66	89.96	80.87	82.18	82.02	83.32	81.04
ZSCL-M - SMILER - large																		
	en	ar	de	es	fa	fr	it	ko	nl	pl	ru	sv	uk	pt	average	ita_avg	nld_avg	fas_avg
mediator	61.78	82.99	86.36	86.70	92.48	85.24	95.61	82.09	95.09	84.06	92.66	94.48	89.34	86.44	86.81	85.23	90.29	85.85
tur	66.99	77.81	83.35	85.74	83.86	83.12	89.33	82.06	88.48	91.88	96.30	79.12	96.25	86.80	85.08	87.61	83.52	81.24
por	60.79	74.95	91.93	86.23	80.32	89.30	89.43	82.18	95.35	86.15	92.75	93.15	86.25	87.23	85.43	84.12	92.43	79.15
avg	57.51	95.40	67.81	61.43	91.71	50.52	58.61	92.75	75.53	70.30	90.38	58.30	88.75	59.42	72.74	69.49	63.04	93.29
random	61.58	81.92	85.29	82.37	85.52	82.00	84.12	86.97	89.41	90.18	95.39	86.31	95.44	82.72	84.94	84.54	85.75	84.80

Figure 11: Detailed transfer performances for SMILER task in ZSCL-M setting.

		DANN - MINION - small										
		english	hindi	japanese	korean	polish	portuguese	spanish	turkey	average	tur_ave	por_ave
tur	mediator	69.50	58.19	39.44	50.36	64.37	63.78	65.08	75.39	60.76	55.84	65.68
	target	60.36	59.67	50.40	54.41	60.75	51.69	61.15	76.16	59.32	60.16	58.49
	avg	77.25	56.52	32.97	52.45	73.85	67.19	64.39	59.73	60.54	50.42	70.67
		DANN - MINION - base										
		english	hindi	japanese	korean	polish	portuguese	spanish	turkey	average	tur_ave	por_ave
tur	mediator	69.51	64.25	51.92	58.64	75.31	69.72	74.07	77.31	67.59	63.03	72.15
	target	65.53	62.93	55.34	59.11	68.86	55.35	67.01	77.18	63.91	63.64	64.19
	avg	80.08	63.11	47.52	54.49	80.38	71.07	73.52	67.93	67.26	58.26	76.26
		DANN - MINION - large										
		english	hindi	japanese	korean	polish	portuguese	spanish	turkey	average	tur_ave	por_ave
tur	mediator	70.72	62.07	46.67	55.89	71.28	65.00	68.88	74.37	64.36	59.75	68.97
	target	68.16	60.20	52.30	56.20	64.55	55.12	67.65	75.86	62.51	61.14	63.87
	avg	79.73	58.97	39.86	53.34	76.90	67.64	71.94	64.13	64.06	54.08	74.05

Figure 12: Detailed transfer performances for MINION task of DANN baseline.

		DANN - SMiLER - small																	
		en	ar	de	es	fa	fr	it	ko	nl	pt	ru	sv	uk	pt	average	ita_ave	nld_ave	fas_ave
tur	mediator	46.19	60.97	63.27	64.73	88.17	61.84	91.30	60.10	93.47	60.83	71.29	61.64	58.71	63.23	67.55	65.18	70.06	69.75
	target	46.98	68.79	68.64	69.24	61.77	67.40	91.53	59.19	77.09	93.34	84.10	63.30	92.17	69.46	72.36	78.12	69.11	63.25
	avg	39.79	48.84	89.89	61.20	49.49	87.64	60.13	59.15	93.44	60.57	65.87	62.46	57.69	60.65	64.06	57.99	83.36	52.49
tur	mediator	35.53	96.44	53.49	45.30	88.57	38.27	42.87	89.56	45.35	49.75	64.93	39.84	62.23	42.93	56.79	49.08	44.24	91.52
	target	42.12	68.76	68.82	60.12	72.00	63.79	71.46	67.00	77.34	66.12	71.55	56.81	67.70	59.07	65.19	62.59	66.69	69.25
	avg	54.03	72.73	76.31	76.74	89.99	71.46	90.02	67.93	92.61	77.97	80.13	75.92	72.69	78.41	76.92	75.71	79.08	76.88
		DANN - SMiLER - base																	
		en	ar	de	es	fa	fr	it	ko	nl	pt	ru	sv	uk	pt	average	ita_ave	nld_ave	fas_ave
tur	mediator	50.52	68.01	74.69	71.34	65.24	70.39	91.94	61.73	80.58	94.20	90.73	65.30	90.43	73.40	74.89	80.37	72.74	64.99
	target	56.85	63.54	90.55	76.48	64.13	88.42	82.30	75.58	94.01	79.49	83.72	77.14	67.25	78.41	76.99	74.93	87.53	67.75
	avg	46.40	95.49	59.61	54.22	89.97	41.76	46.75	91.64	63.16	64.15	74.77	48.35	65.58	50.92	63.77	57.54	53.22	92.37
tur	mediator	51.95	74.94	75.29	69.69	77.33	68.01	77.75	74.22	82.59	78.95	82.34	66.68	73.99	70.28	73.14	72.14	73.14	75.50
	target	58.30	79.47	84.74	82.95	92.01	82.25	93.11	79.22	94.10	83.07	89.96	80.16	84.02	84.95	83.45	82.34	85.31	83.57
	avg	55.95	77.77	81.58	80.54	78.25	81.31	93.60	76.54	87.60	95.15	94.39	74.20	94.26	83.50	82.47	85.34	81.17	77.52
		DANN - SMiLER - large																	
		en	ar	de	es	fa	fr	it	ko	nl	pt	ru	sv	uk	pt	average	ita_ave	nld_ave	fas_ave
tur	mediator	58.32	72.27	91.08	83.06	72.19	88.39	86.78	81.04	94.12	84.23	90.17	81.60	81.59	84.45	82.09	81.23	88.80	75.17
	target	47.11	95.49	62.02	54.89	92.59	46.69	51.83	92.21	62.07	62.39	80.91	49.92	75.07	53.51	66.19	60.82	55.17	93.43
	avg	54.92	81.25	79.85	75.36	83.76	74.66	81.33	82.25	84.47	81.21	88.86	71.47	83.73	76.60	78.55	77.43	77.61	82.42

Figure 13: Detailed transfer performances for SMiLER task of DANN baseline.

ZSCL-R - MINION - small											
	english	hindi	japanese	korean	polish	portuguese	spanish	turkey	average	tur_ave	por_ave
mediator	72.15	62.22	49.23	56.95	67.60	67.54	68.43	77.44	65.19	61.46	68.93
tur	66.81	60.83	56.53	59.45	64.68	56.01	64.84	76.50	63.21	63.33	63.09
por	79.00	61.64	42.51	56.35	75.98	68.21	69.76	65.69	64.89	56.55	73.24

ZSCL-R - MINION - base											
	english	hindi	japanese	korean	polish	portuguese	spanish	turkey	average	tur_ave	por_ave
mediator	76.19	66.67	58.29	61.79	77.05	71.57	75.85	79.18	70.82	66.48	75.16
tur	72.83	65.30	59.76	62.07	73.28	61.02	71.31	78.62	68.03	66.44	69.61
por	82.18	66.48	50.25	60.52	81.85	72.74	76.09	71.73	70.23	62.25	78.21

ZSCL-R - MINION - large											
	english	hindi	japanese	korean	polish	portuguese	spanish	turkey	average	tur_ave	por_ave
mediator	76.22	64.85	51.34	59.72	74.82	69.83	73.24	77.16	68.40	63.27	73.53
tur	72.32	63.28	59.36	59.84	70.17	58.30	69.63	77.25	66.27	64.93	67.61
por	81.23	63.09	48.84	58.83	80.60	71.66	74.42	67.31	68.25	59.52	76.98

Figure 14: Detailed transfer performances for MINION task in ZSCL-R setting.

ZSCL-R - SMiLER - small																		
	en	ar	de	es	fa	fr	it	ko	nl	pt	ru	sv	uk	pt	average	ita_ave	nld_ave	fas_ave
mediator	56.56	72.75	82.15	81.73	89.97	79.60	95.10	73.74	94.70	82.27	85.12	82.49	78.92	81.47	81.18	80.17	84.73	78.82
it	56.07	80.25	80.05	81.10	69.14	80.06	95.81	69.60	87.61	94.73	95.05	74.93	94.26	81.89	81.47	85.56	80.66	73.00
nl	59.82	67.26	91.13	83.59	68.20	89.01	85.89	78.13	94.91	82.73	86.32	91.60	74.47	81.90	81.07	79.25	91.66	71.20
fr	40.84	96.46	54.53	47.98	90.29	40.73	49.19	91.38	66.95	53.23	82.35	45.20	67.34	45.01	62.25	55.13	51.85	92.71
avg	53.32	79.18	76.97	73.60	79.40	72.35	81.50	78.21	86.04	78.24	87.21	73.56	78.75	72.57	76.49	75.03	77.23	78.93

ZSCL-R - SMiLER - base																		
	en	ar	de	es	fa	fr	it	ko	nl	pt	ru	sv	uk	pt	average	ita_ave	nld_ave	fas_ave
mediator	61.13	82.17	85.15	83.25	94.17	81.49	96.03	77.87	95.30	82.33	85.38	91.87	79.15	84.62	84.28	81.70	88.45	84.74
it	57.23	75.92	82.52	80.76	75.70	79.98	96.06	72.05	87.22	95.40	94.57	77.31	96.25	82.76	82.41	86.15	81.76	74.56
nl	59.88	71.19	91.62	84.74	71.36	89.38	87.43	79.94	95.32	82.86	89.49	94.13	82.54	83.85	83.12	81.54	92.61	74.16
fr	53.68	96.07	66.76	57.22	93.42	47.53	52.85	93.02	73.62	66.70	84.69	55.80	84.30	54.59	70.02	64.86	60.93	94.17
avg	57.98	81.34	81.51	76.49	83.66	74.59	83.09	80.72	87.86	81.82	88.53	79.78	85.56	76.45	79.96	78.56	80.94	81.91

ZSCL-R - SMiLER - large																		
	en	ar	de	es	fa	fr	it	ko	nl	pt	ru	sv	uk	pt	average	ita_ave	nld_ave	fas_ave
mediator	60.58	77.87	87.16	86.54	93.98	83.30	95.61	79.65	94.86	84.55	91.50	91.22	88.26	86.07	85.80	84.73	89.13	83.83
it	57.95	80.49	83.92	83.11	88.18	83.15	95.37	80.58	88.96	95.72	96.66	81.46	96.25	85.33	85.51	87.20	84.37	83.08
nl	60.93	76.14	92.57	85.18	77.81	89.50	88.83	83.25	95.59	84.94	93.41	93.48	86.98	86.22	85.34	83.78	92.78	79.07
fr	59.54	96.92	70.94	64.71	94.64	51.81	61.94	94.00	76.89	72.22	90.67	63.64	89.54	61.00	74.89	71.37	65.82	95.19
avg	59.75	82.86	83.65	79.89	88.65	76.94	85.44	84.37	89.07	84.36	93.06	82.45	90.26	79.66	82.89	81.77	83.03	85.29

Figure 15: Detailed transfer performances for SMiLER task in ZSCL-R setting.