

Persian Natural Language Inference: A Meta-learning approach

Anonymous ACL submission

Abstract

Incorporating information from other languages can improve the results of tasks in low-resource languages. A powerful method of building functional natural language processing systems for low-resource languages is to combine multilingual pre-trained representations with cross-lingual transfer learning. In general, however, shared representations are learned separately, either across tasks or across languages. This paper proposes a meta-learning approach for inferring natural language in Persian. Alternately, meta-learning uses different task information (such as QA in Persian) or other language information (such as natural language inference in English). Also, we investigate the role of task augmentation strategy for forming additional high-quality tasks. We evaluate the proposed method using four languages and an auxiliary task. Compared to the baseline approach, the proposed model consistently outperforms it, improving accuracy by roughly six percent. We also examine the effect of finding appropriate initial parameters using zero-shot evaluation and CCA similarity.

1 Introduction

In natural language processing (NLP), the goal is to improve models for the processing and production of human languages. As part of NLP, several tasks are defined, each of which covers a different level of natural language understanding. Meanwhile, natural language inference (NLI) is considered an appropriate and rigorous measure of language comprehension. This task requires to recognize the consequences of natural language sentences, which indicates how well it understands the language (MacCartney, 2009).

NLI aims to determine the inferential relationship between a premise p and a hypothesis h . The problem involves a three-class classification in which every pair (p, h) falls into one of three categories: entailment, contradiction, and neutral. If

the hypothesis can be inferred from the premise, pair (p, h) will be assigned to the entailment class. For a hypothesis that contradicts the premise, pair (p, h) will be assigned to the contradiction and neutral otherwise (Amirkhani et al., 2020).

As we know, the Persian language lacks sufficient linguistic resources when it comes to natural language understanding. The lack of data can be addressed by collecting annotated data, but this process is both time-consuming and expensive (Nooralahzadeh et al., 2020). FarsTail (Amirkhani et al., 2020) is currently available for Persian, which is created using the same method as SciTail (Khot et al., 2018). It contains 10,367 samples. Also, ParsiNLU (Khashabi et al., 2021) is created for high-level tasks in Persian and for NLI, it consists of 2700 samples. As it turns out, this amount of data is too small compared with resource-rich languages (such as English, which has only 550,000 samples in the SNLI (Bowman et al., 2015) dataset).

Researchers have tried to solve the data scarcity problem by using cross-language methods. Recent work on cross-lingual learning has mainly focused on transfer between languages already covered by pre-trained representations (Wu and Dredze, 2019). Nonetheless, these techniques do not readily transfer to low-resource languages in which (1) large monolingual corpora are unavailable for pre-training, and (2) sufficient labeled data is lacking for fine-tuning downstream tasks (Xia et al., 2021).

The results of experimental studies for Persian using different embedding methods including word2vec (Mikolov et al., 2013), fastText (Bojanowski et al., 2017), ELMo (Peters et al., 2018), and BERT (Devlin et al., 2019) and various models, such as, DecompAtt (Parikh et al., 2016), ESIM (Chen et al., 2016), HBMP (Talman et al., 2019), and ULMFiT (Howard and Ruder, 2018) is reported in FarsTail (Amirkhani et al., 2020). Although this cross-lingual information sharing has

182 using monolingual data and the other using paral- 232
183 lel data and a new cross-lingual language model 233
184 objective. Singh et al. (2019) introduced a cross- 234
185 lingual data augmentation method that substitutes 235
186 part of the input text with its translation in another 236
187 language. 237

188 Huertas-Tato et al. (2021) designed a new archi- 238
189 tecture called Siamese Inter-Lingual Transformer 239
190 (SILT), to align multilingual embeddings for NLI 240
191 efficiently. The paper points out that transformer 241
192 models are unable to generalize to other domains 242
193 and have problems with multilingual and inter- 243
194 linguistic scenarios. A new network has been de- 244
195 veloped to overcome these weaknesses by combin- 245
196 ing three parts: a multilingual transformer as pre- 246
197 trained embedding, an alignment matrix to com- 247
198 pute the similarity between two sentences, and a 248
199 multi-head self-attention block to interpret input 249
200 strings. 250

201 Despite the advances that Cross-lingual methods 251
202 have made, building NLP systems in these settings 252
203 is challenging for several reasons. First, the tar- 253
204 get language does not contain sufficient annotated 254
205 data for effective fine-tuning. Secondly, pre-trained 255
206 multilingual representations are not directly trans- 256
207 ferable due to language disparities (Xia et al., 2021). 257
208 In contrast to these methods, we consider setting 258
209 up training models simultaneously on multiple lan- 259
210 guages and tasks. 260

211 2.2 Meta-learning 262

212 Meta-learning addresses the problem of learning 263
213 to learn. By examining many learning problems, 264
214 a meta-learner learns a model (Liu et al., 2020). 265
215 Specifically, the meta-learner uses a meta training 266
216 set $\text{MS} = \{(\mathbb{S}_i^s, \mathbb{T}_i^s)\}_{i=1}^{N^s}$, where $(\mathbb{S}_i^s, \mathbb{T}_i^s)$ are the 267
217 training (support) and test (query) set of the i^{th} 268
218 learning problem and N^s is the number of learn- 269
219 ing problems used for training; and a meta test set 270
220 $\text{MT} = \{(\mathbb{S}_i^t, \mathbb{T}_i^t)\}_{i=1}^{N^t}$, where $(\mathbb{S}_i^t, \mathbb{T}_i^t)$ are the sup- 271
221 port and query set of the i^{th} test learning problem, 272
222 while N^t is the number of learning problems used 273
223 for the test. Given MS , the meta-learner learns 274
224 how to map a pair (\mathbb{S}, \mathbb{T}) into an algorithm that 275
225 leverages \mathbb{S} to optimally solve \mathbb{T} . 276

226 Due to the lack of well-defined task distribution, 277
227 meta-learning has not yet succeeded in NLP, lead- 278
228 ing to attempts that treat datasets as tasks. An ad 279
229 hoc task distribution causes problems with quantity 280
230 and quality. Murty et al. (2021) provide a way to 281
231 break down heterogeneous tasks such as datasets

232 into a set of appropriate subtasks. With this method, 233
234 data is transferred to the feature space using a pre- 235
236 trained model. They use k-means to decompose 236
237 data into k clusters and create tasks by combining 238
239 these clusters. 240

241 Recently, however, the combination of cross- 242
243 lingual techniques in meta-learning frameworks 243
244 has also been extensively studied. To train a model 244
245 for low-resource languages on NLI and QA tasks, 245
246 Nooralahzadeh et al. (2020) uses the MAML al- 246
247 gorithm and auxiliary languages. van der Heij- 247
248 den et al. (2021) study the text documents classi- 248
249 fication problem in monolingual and multilingual 249
250 modes, using different algorithms such as, Pro- 250
251 totypical Networks (Snell et al., 2017), MAML 251
(Finn et al., 2017), Reptile (Nichol et al., 2018), 252
253 and ProtoMAML (Triantafillou et al., 2019). Also, 253
254 Tarunesh et al. (2021) examine the interaction be- 254
255 tween different languages and tasks to learn an 255
256 appropriate common feature space. 256

257 Additionally, transfer-learning can be helpful for 257
258 low-resource languages. Xia et al. (2021) introduce 258
259 a meta-learning-based framework called MetaXL 259
260 for extremely low-resource languages. MetaXL 260
261 learns an intelligent representational conversion 261
262 from several auxiliary languages to the target lan- 262
263 guage, bringing the feature space of these lan- 263
264 guages closer together for more efficient conver- 264
265 sion. The main idea is to use a Representation 265
266 Transformation network between the main model 266
267 layers which are trained only with target language 267
268 data. 268

269 To the best of our knowledge, this paper is the 269
270 first attempt to study meta-learning for solving the 270
271 NLI problem in the Persian language. Also, we 271
272 are pioneers in using task-language pairs as meta- 272
273 learning tasks in the Persian language. 273

274 2.3 Task Augmentation 275

276 Machine learning algorithms usually assume that 276
277 the train and test data have the same distribution. In 277
278 contrast, the meta-learning framework treats tasks 278
279 as training examples and trains a model to adapt to 279
280 all of them. Meta-learning also assumes that the 280
281 training and new tasks are drawn from the same 281
282 distribution of tasks $p(\tau)$. In NLP, datasets are 282
283 typically treated as tasks, and meta-learners are 283
284 then overfitting their adaptation mechanisms. NLP 284
285 datasets are highly heterogeneous, which causes 285
286 many learning episodes to have the poor transfer 286
287 between their support and query sets, which dis- 287
288 289
290
291

suades meta-learners from adapting (Murty et al., 2021).

To deal with overfitting challenges, Yin et al. (2019) propose a meta-regularizer to mitigate memorization overfitting, but don’t study learner overfitting. Rajendran et al. (2020) study task augmentation for mitigating meta-learners overfitting in the context of few-shot label adaptation. SMLMT method (Bansal et al., 2020) creates new self-supervised tasks that improve meta-overfitting, but this does not directly address the dataset-as-tasks problem. In contrast, the DReCa method (Murty et al., 2021) addresses the dataset-as-tasks problem and focuses on using clustering as a way to subdivide and fix tasks that already exist. In this paper, we use DReCa as a task augmentation strategy for our method since it mitigates meta-overfitting without any additional unlabeled data.

	NLI	QA
FA	FarsTail (10.3K)	PersianQA (9K)
EN	XNLI (392k)	—
ES	tr. XNLI (392k)	—
DE	tr. XNLI (392k)	—
FR	tr. XNLI (392k)	—

Table 1: Overview of datasets from a variety of sources. For the NLI task, we use the XNLI dataset for English, and their translated versions (tr.) for Spanish(ES), German(DE), and French(FR) provided in XTREME. For each dataset, the number of training instances is also mentioned.

3 The Proposed Methodology

In our setting, firstly, we prepare a set of task-language pairs to provide meta-learning tasks. Afterward, in each episode, we sample some tasks and feed them to the meta-learner. In the rest of this section, we describe the proposed task sampling strategy, the proposed meta-learning algorithm, and the proposed task augmentation strategy.

3.1 The Proposed Task Sampling Strategy

In meta-learning, task selection has a profound impact on model performance. For this reason, we create a queue of tasks first. We can create this queue using different scenarios such as selecting languages for a target task (Gu et al., 2018), selecting tasks for a target language (Dou et al., 2019), and picking from various auxiliary languages and auxiliary tasks. In the meta-training section, we sample some tasks from the queue. Formally, the queue’s

tasks are represented by \mathcal{D} . We need to sample tasks from \mathcal{M} , which is a Multinomial distribution over $P_{\mathcal{D}}(i)$ s. Thus, we investigate temperature-based heuristic sampling (Aharoni et al., 2019), which defines the probability of any dataset as a function of its size as,

$$P_{\mathcal{D}}(i) = q_i^{1/\tau} / \left(\sum_{k=1}^n q_k^{1/\tau} \right) \quad (1)$$

where $P_{\mathcal{D}}(i)$ is the probability of sampling the i th task, q_i is the size of i th task, and τ is the temperature parameter. With $\tau = 1$, tasks are randomly sampled proportionately to their dataset sizes, and with $\tau \rightarrow \infty$, they follow a uniform distribution.

3.2 The Proposed Meta-learning Algorithms

Meta-learning is the process of building a model that can solve a new task with only a few labeled examples by training on a variety of tasks with rich annotations. The key idea is to train the model’s initial parameters such that the model has maximal performance on a new task after the parameters have been updated through zero or a couple of gradient steps (Yin, 2020). MAML (model-agnostic meta-learning) (Finn et al., 2017) is one of the most significant algorithms. We describe one episode of the MAML algorithm in Appendix A.1. MAML is quite difficult to train, since there are two levels of training. Therefore, we use the following two optimization-based and metric-based meta-learning algorithms in this work.

Reptile (Nichol et al., 2018) is a first-order optimization-based algorithm that moves weights toward a manifold of the weighted averages of task-specific parameters $\theta_i^{(m)}$. It samples training tasks from $p(\mathcal{T}) : \tau_1, \dots, \tau_i, \dots, \tau_n$. For each training task, it generates an episode that just contains the support set data. For training task τ_i , let’s assume the original parameters θ have gone through m steps of updating and become $\theta_i^{(m)}$ (i.e., $\theta_i^{(m)} = \text{AdamW}(L_{\tau_i}, \theta, m)$ (2)), then Reptile updates θ as follows (Yin, 2020):

$$\theta \leftarrow \theta + \beta \frac{1}{|\{\mathcal{T}\}|} \sum_{\tau_i \sim \mathcal{M}} \left(\theta_i^{(m)} - \theta \right) \quad (3)$$

Prototypical Networks (Snell et al., 2017) is a metric-based meta-learning algorithm. Prototypical networks learn a metric space in which classification can be performed by computing distances to prototype representations of each class. In general,

they are composed of an embedding network f_θ and a distance function $d(x_1, x_2)$. Using the following equation, the embedding network encodes the support set samples S_c and computes prototypes μ_c per class based on the mean sample encodings for that class.

$$\mu_c = \frac{1}{|S_c|} \sum_{(x_i, y_i) \in S_c} f_\theta(x_i). \quad (4)$$

A Prototypical network classifies a new sample according to the following rule.

$$p(y = c | x) = \frac{\exp(-d(f_\theta(x), \mu_c))}{\sum_{c' \in C} \exp(-d(f_\theta(x), \mu_{c'}))} \quad (5)$$

Thus, we define the *distance-based cross entropy* (DCE) loss as follows:

$$\text{Loss}(DCE) = -\log P(y = c | x) \quad (6)$$

To ensure that the feature space is robust to noise, we also use the Cross Entropy (CE) loss (more details can be found in Appendix A.3.1).

3.3 The Task Augmentation Strategy

First, we use dataset-as-tasks strategy that is the most common method for selecting tasks for meta-learning in NLP applications. Next, we employ DReCa to form additional high quality tasks. The goal of DReCa is to take a heterogeneous task (such as a dataset) and produce a decomposed set of tasks. Given a training task T_i^{tr} , DReCa first groups examples by their labels, and then embeds examples within each group with an embedding function $\text{EMBED}(\cdot)$. Concretely, for each N -way classification task T_i^{tr} , it forms groups $g_l^i = \{(\text{EMBED}(x_i), y_i) \mid y_i = l\}$. Then, it proceeds to refine each label group into K clusters via k-means clustering to break down T_i^{tr} into groups $\{C^j(g_l^i)\}_{j=1}^K$ for $l = 1, 2, \dots, N$. These cluster groups can be used to produce K^N potential DReCa tasks. Each task is obtained by choosing one of K clusters for each of the N label groups, and taking their union.

4 Experimental Setup

4.1 Datasets

We use FarsTail (Amirkhani et al., 2020) for the target dataset. FarsTail is the only large-scale Persian corpus for the NLI task, with 10,367 samples. The samples are generated from 3,539 multiple-choice

questions with the least amount of annotators' interventions or selected from natural sentences that already exist independently in the wild, similarly to the SciTail dataset (Khot et al., 2018).

We also use XTREME (Hu et al., 2020) as an auxiliary dataset. XTREME is a multilingual multi-task benchmark consisting of classification, structured prediction, QA, and retrieval tasks. We use this benchmark to prepare NLI data for auxiliary languages. Note that, large-scale datasets for NLI were only available in English. However, the authors of XTREME developed a custom-built translation system to get translated datasets for NLI. Furthermore, we consider the QA as an auxiliary task. Therefore, we use PersianQA (Ayoubi and Davoodeh, 2021) which is a Persian reading comprehension dataset for QA, containing over 9000 entries. Table 1 summarizes the employed dataset specifications.

4.2 Baselines

On the FarsTail dataset, Amirkhani et al. (2020) present results of various traditional and deep learning-based methods. According to the results of this paper, the highest test accuracy is obtained by using a translation-based approach, i.e., *Translate-Source* with fastText embeddings. In *Translate-Source*, the Persian-translated MultiNLI training set is combined with FarsTail training data for training an ESIM model. Furthermore, FarsTail's authors reported mBERT fine-tuning results in FarsTail webpage¹. Therefore, we use these results as baselines.

4.3 Implementation Details

In this study, we aim to compare the effects of meta-learning algorithms on classification accuracy with those of fine-tuning and non-episodic algorithms. To make a fair comparison, we first fine-tune our pre-trained models using training data of the auxiliary task in a non-episodic approach. Afterward, we fine-tune the obtained model using the training data of the target task. In this approach, we use mBERT (Devlin et al., 2019), and XLM-R (Conneau et al., 2020), which are known as the state-of-the-art multilingual pre-trained models, and ParsBERT (Farahani et al., 2021) as a monolingual transformer-based model for the Persian language.

In the meta-learning approach, we use the XLM-R model with output layers tailored for each task

¹<https://github.com/dml-gom/FarsTail>

and train it with Reptile and Prototypical algorithms. To select the hyperparameters of the Reptile algorithm, we utilize the experiments done in Tarunesh et al. (2021). Appendix A.2 provides further details. The Prototypical algorithm is used only in cross-lingual experiments, and we use Euclidean distance as its distance function. The auxiliary languages are arranged in two scenarios. In the first scenario, support and query set data are generated from auxiliary languages, while in the second scenario, the query set is drawn from both auxiliary and target languages. Detailed information is provided in Appendix A.3.2.

Furthermore, we fine-tune the obtained models on Persian training data using the following two methods. The first method is **non-episodic**, which involves fine-tuning models in batches. The second method is **episodic**, in which episodes are constructed first, and then the models are fine-tuned according to the algorithm used.

5 Results

The meta-learning model is tested on different combinations and configurations of the auxiliary tasks. The accuracy results of the Reptile algorithm are presented in Table 2. In addition to the zero-shot and fine-tuning results, we report the accuracy of another scenario. In this scenario, training data of the target language is placed in the meta-training stage along with other auxiliary tasks and cooperate in a training process. Consequently, this scenario does not involve fine-tuning phase. The results of the mentioned scenario are shown in the last column of Table 2.

Table 3 shows the accuracy scores using the Prototypical Network. In the first section of this table, we generate both support and query sets from Persian language data, without using auxiliary tasks. In the second section of this table, the results of the first multi-lingual scenario (where both the support and query sets are generated from auxiliary languages data) are reported in rows 5 to 12. In rows 13 to 16, we show the results of the second multi-lingual scenario (where the support set is drawn from auxiliary language data and the query set is drawn from both auxiliary and Persian language data). Lastly, we added the DReCa strategy and presented the results in rows 17 to 20.

Additionally, we conducted zero-shot evaluations of both algorithms. Zero-shot results are presented in the first accuracy column of Tables 2 and

3. The confusion matrices of the best-performing models for both Reptile and Prototypical algorithms are also depicted in Appendix A.4.

6 Discussion and Analysis

Table 2 shows that the multi-lingual models are always better than the multi-task models. Due to the fact that tasks like NLI (which require deeper semantic representations) are more likely to benefit from combining data from different languages. We found that our meta-learned models perform better than baselines and non-episodic models. The reason is that the goal of standard meta-learning is to find a model that generalizes well to a new target task. In addition, we compared two different meta-learning algorithms to evaluate their superiority in this paper. From Tables 2 and 3, we can see that Prototypical performed better than Reptile. It is because Prototypical networks use class representations instead of example representations. Therefore, it finds a suitable representation for each class during the meta-train stage.

As part of another experiment, we combined data from the target language with data from other auxiliary tasks for meta-training. Based on the results of these experiments (last column of Table 2 for Reptile and rows 13 to 16 of Table 3 for Prototypical), the model’s accuracy has decreased. This is due to the fact that target language data is so small when compared with auxiliary language data. So, unbalanced training data confuses the training process and decreases the model’s accuracy. In any case, the cooperation of the target language during the training process is a great idea for future work.

As indicated in the last two columns of Table 3, episodic fine-tuning is significantly superior to normal fine-tuning. It demonstrates that episodic training is effective even on single language data and creates a generality in the level of training and test data.

We examined the proximity between the feature spaces of the auxiliary languages and the target language quantitatively and qualitatively. At first, we collect representations of the auxiliary and target languages from non-episodic, Reptile, and Prototypical models. In Fig. 1, we present 2-component PCA visualization for comparison. We also evaluated the models using a distance metric commonly used in vision and NLP tasks (Huttenlocher et al., 1993; Dubuisson and Jain, 1994; Patra et al., 2019; Xia et al., 2021). Informally, the

Row	Model	Shot	Aux. Tasks	Zero-shot	non-episodic fine-tune	Add NLI-fa in m.t.
Baselines						
1	Translate-Source*	—	—	—	78.13	—
2	mBERT*	—	—	—	83.38	—
Non-episodic approach						
3	ParsBERT	—	—	—	74.64	—
4	mBERT	—	—	—	81.95	—
5		—	nli-en	56.53	81.38	—
6		—	nli-(en, es, de, fr)	67.88	82.34	—
7	XLM-R	—	—	—	81.97	—
8		—	nli-en	69.49	86.55	—
9		—	nli-(en, es, de, fr)	69.09	84.69	—
Meta-learning approach						
10	XLM-R	1	—	64.19	84.31	83.37
11		4	—	70.96	87.17	86.00
12		8	nli-en	70.70	87.11	86.65
13		16	—	71.03	87.43	86.52
14		1	—	65.17	85.21	83.91
15		4	nli-(en, es, de, fr)	72.27	87.57	85.74
16		8	—	71.61	88.35	88.22
17		16	—	71.22	88.02	87.76
18		1	—	34.18	81.48	81.58
19		4	qa-fa	34.18	81.38	84.96
20		8	—	33.79	82.14	83.59
21		16	—	34.18	82.23	84.70
22		1	—	46.42	83.53	85.16
23		4	nli-en,	66.02	86.52	86.26
24		8	qa-fa	64.26	86.98	86.46
25		16	—	46.88	86.52	86.13

Table 2: Average test accuracy of the Reptile algorithm with baselines and non-episodic approach results on the Persian NLI task. The first accuracy column shows results before fine-tuning on the Persian NLI train-set (called zero-shot). In the second accuracy column, we provided results after fine-tuning on the Persian NLI train-sets. The last accuracy column reports results of using the Persian NLI train-set in the meta-training phase (m.t.). The data with * comes from FarsTail’s paper and webpage.

Row	Model	Shot	Support	Query	Zero-shot	non-episodic fine-tune	episodic fine-tune	
Without auxiliary tasks								
1	XLM-R	1	—	—	—	70.38	79.30	
2		4	nli-fa	nli-fa	—	81.97	85.22	
3		8	—	—	—	83.98	84.64	
4		16	—	—	—	85.29	85.74	
With auxiliary tasks								
5	XLM-R	1	—	—	68.10	84.83	86.07	
6		4	nli-en	nli-en	70.57	86.72	87.50	
7		8	—	—	70.77	86.72	87.37	
8		16	—	—	73.18	87.76	88.54	
9		1	—	—	69.15	85.01	85.97	
10		4	nli-(en, es, de, fr)	nli-(en, es, de, fr)	70.25	86.78	87.63	
11		8	—	—	71.09	88.48	89.39	
12		16	—	—	72.20	88.15	88.28	
13		1	—	—	—	84.15	85.12	
14		4	nli-(en, es, de, fr)	nli-(en, es, de, fr, fa)	—	86.33	86.78	
15		8	—	—	—	86.33	86.46	
16		16	—	—	—	86.78	87.24	
17		XLM-R+ DReCa	8	nli-en	nli-en	70.44	87.96	88.87
18			16	—	—	71.94	87.24	88.74
19			8	nli-(en, es, de, fr)	nli-(en, es, de, fr)	71.16	87.74	88.48
20			16	—	—	71.61	87.30	88.22

Table 3: Average test accuracy on the Persian NLI task using Prototypical algorithm with and without auxiliary tasks. The last accuracy column reports results after episodic fine-tuning on the Persian NLI train-set.

Hausdorff distance measures the distance between data representations of auxiliary languages and the target language. Given a set of representations of the auxiliary language $\mathcal{S} = \{s_1, s_2, \dots, s_m\}$ and a set of representations of the target language $\mathcal{T} = \{t_1, t_2, \dots, t_m\}$ we compute the Hausdorff

distance as follows:

$$\max \left\{ \max_{s \in \mathcal{S}} \min_{t \in \mathcal{T}} d(s, t), \max_{t \in \mathcal{T}} \min_{s \in \mathcal{S}} d(s, t) \right\} \quad (7)$$

where cosine distance is used as the inner distance, i.e.,

$$d(s, t) \triangleq 1 - \cos(s, t) \quad (8)$$

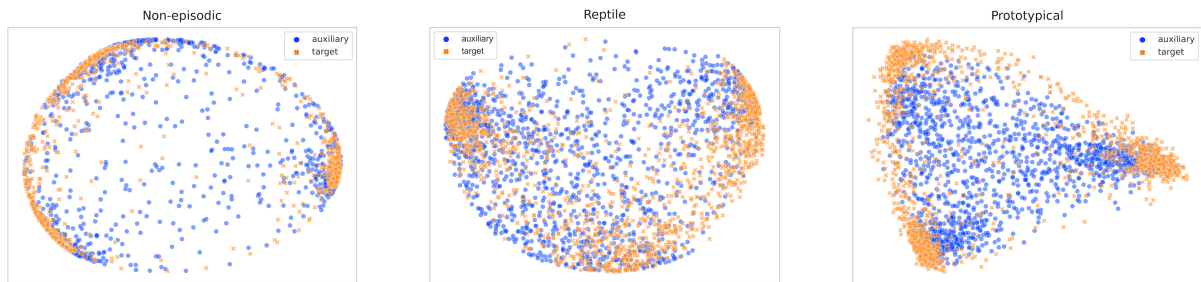


Figure 1: PCA visualization of non-episodic, Reptile, and Prototypical models to examine the closeness of the auxiliary and target languages feature spaces.

Compared to the non-episodic method, we observe a drastic drop of Hausdorff distance from 0.18 to 0.05 for Prototypical and also, we see a minor decline of Hausdorff distance from 0.18 to 0.13 for Reptile. Both qualitative visualization and quantitative metrics confirm that meta-learning approaches bring the distributions of auxiliary and target language data closer together, thus increasing the accuracy on the target language.

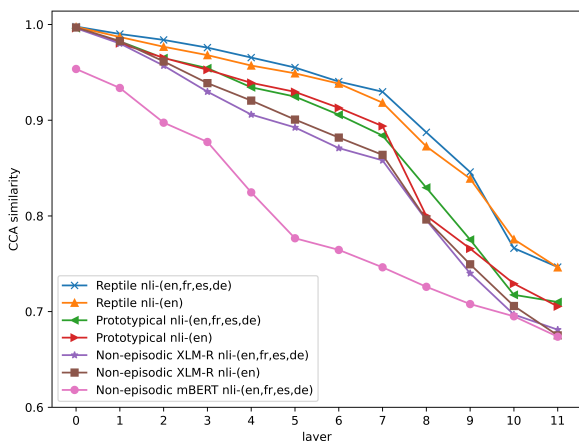


Figure 2: CCA similarity for each transformer layer. We calculate the similarity before and after fine-tuning on the FarsTail training data.

The advantage of meta-learning methods is that they obtain the appropriate initial parameters for the target language, as mentioned. The zero-shot test is used as a criterion to evaluate this point, and it shows that meta-learning-based models are more accurate than other methods. The generality of the initial parameters can also be assessed via canonical correlation analysis (CCA) (Raghu et al., 2017; Morcos et al., 2018). Using this criterion, we compare the output of each layer before and after fine-tuning, and the results are presented in Fig. 2. The meta-learning models have a higher

CCA similarity, which indicates the model obtained more general parameters before fine-tuning.

In the next experiment, we apply the DReCa strategy and train the model with the Prototypical algorithm. According to Table 3, some results have improved, while others have remained the same. It illustrates that task augmentation in meta-algorithms affects the model’s accuracy. However, defining the appropriate task augmentation strategy still needs research.

7 Conclusion

We present effective use of meta-learning to benefit from other tasks or languages. We advantageously leverage this approach to improve NLI in Persian as a low-source language. We found that our meta-learning model outperformed competitive baseline models. In response to the concept of treating entire datasets as tasks, we use DReCa as a general-purpose task augmenting approach. Finally, zero-shot evaluations illustrate the generality of the results obtained by meta-learning. This work will be extended to other cross-lingual NLP tasks in Persian in the future. Furthermore, we would like to use a self-supervised approach to provide a useful starting point for parameters.

References

- Roe Aharoni, Melvin Johnson, and Orhan Firat. 2019. [Massively multilingual neural machine translation](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics.
- Hossein Amirkhani, Mohammad Azari Jafari, Azadeh Amirak, Zohreh Pourjafari, Soroush Faridan Jahromi, and Zeinab Kouhkan. 2020. [FarsTail: A persian natural language inference dataset](#). *arXiv preprint arXiv:2009.08820*.

622	Sajjad Ayoubi and Mohammad Yasin Davoodeh. 2021.	Mehrdad Farahani, Mohammad Gharachorloo,	677
623	Persianqa: a dataset for persian question answering.	Marzieh Farahani, and Mohammad Manthouri.	678
624	Trapit Bansal, Rishikesh Jha, Tsendsuren Munkhdalai,	2021. ParsBERT: Transformer-based model for	679
625	and Andrew McCallum. 2020. Self-supervised	persian language understanding. <i>Neural Process.</i>	680
626	meta-learning for few-shot natural language classifica-	<i>Lett.</i>	681
627	tion tasks. In <i>Proceedings of the 2020 Conference</i>	Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017.	682
628	<i>on Empirical Methods in Natural Language Process-</i>	Model-agnostic meta-learning for fast adaptation of	683
629	<i>ing (EMNLP).</i> Association for Computational Lin-	deep networks. In <i>International Conference on Ma-</i>	684
630	guistics.	<i>chine Learning.</i> PMLR.	685
631	Piotr Bojanowski, Edouard Grave, Armand Joulin, and	Jiatao Gu, Yong Wang, Yun Chen, Victor O. K. Li,	686
632	Tomas Mikolov. 2017. Enriching word vectors with	and Kyunghyun Cho. 2018. Meta-learning for low-	687
633	subword information. <i>Transactions of the Associa-</i>	resource neural machine translation. In <i>Proceedings</i>	688
634	<i>tion for Computational Linguistics.</i>	<i>of the 2018 Conference on Empirical Methods in</i>	689
635	Samuel R. Bowman, Gabor Angeli, Christopher Potts,	<i>Natural Language Processing.</i> Association for Com-	690
636	and Christopher D. Manning. 2015. A large anno-	<i>putational Linguistics.</i>	691
637	tated corpus for learning natural language inference.	Jeremy Howard and Sebastian Ruder. 2018. Universal	692
638	In <i>Proceedings of the 2015 Conference on Empirical</i>	language model fine-tuning for text classification. In	693
639	<i>Methods in Natural Language Processing.</i> Associa-	<i>Proceedings of the 56th Annual Meeting of the Asso-</i>	694
640	tion for Computational Linguistics.	<i>ciation for Computational Linguistics.</i> Association	695
641	Qian Chen, Xiaodan Zhu, Zhenhua Ling, Si Wei, Hui	for Computational Linguistics.	696
642	Jiang, and Diana Inkpen. 2016. Enhanced LSTM	Junjie Hu, Sebastian Ruder, Aditya Siddhant, Gra-	697
643	for natural language inference. <i>arXiv preprint</i>	ham Neubig, Orhan Firat, and Melvin Johnson.	698
644	<i>arXiv:1609.06038.</i>	2020. XTREME: A massively multilingual multi-	699
645	Alexis Conneau, Kartikay Khandelwal, Naman Goyal,	task benchmark for evaluating cross-lingual gener-	700
646	Vishrav Chaudhary, Guillaume Wenzek, Francisco	alization. In <i>International Conference on Machine</i>	701
647	Guzmán, Edouard Grave, Myle Ott, Luke Zettle-	<i>Learning.</i> PMLR.	702
648	moyer, and Veselin Stoyanov. 2020. Unsupervised	Javier Huertas-Tato, Alejandro Martín, and David	703
649	cross-lingual representation learning at scale. In	Camacho. 2021. SML: a new semantic em-	704
650	<i>Proceedings of the 58th Annual Meeting of the Asso-</i>	bedding alignment transformer for efficient cross-	705
651	<i>ciation for Computational Linguistics.</i> Association	lingual natural language inference. <i>arXiv preprint</i>	706
652	for Computational Linguistics.	<i>arXiv:2103.09635.</i>	707
653	Alexis Conneau and Guillaume Lample. 2019. Cross-	Daniel P Huttenlocher, Gregory A. Klanderman, and	708
654	lingual language model pretraining. <i>Advances in</i>	William J Rucklidge. 1993. Comparing images us-	709
655	<i>Neural Information Processing Systems.</i>	ing the hausdorff distance. <i>IEEE Transactions on</i>	710
656	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	<i>pattern analysis and machine intelligence.</i>	711
657	Kristina Toutanova. 2019. BERT: Pre-training of	Daniel Khashabi, Arman Cohan, Siamak Shakeri,	712
658	deep bidirectional transformers for language under-	Pedram Hosseini, Pouya Pezeshkpour, Malihe	713
659	standing. In <i>Proceedings of the 2019 Conference of</i>	Alikhani, Moin Aminnaseri, Marzieh Bitaab, Faeze	714
660	<i>the North American Chapter of the Association for</i>	Brahman, Sarik Ghazarian, Mozdeh Gheini,	715
661	<i>Computational Linguistics: Human Language Tech-</i>	Arman Kabiri, Rabeeh Karimi Mahabagdi, Omid	716
662	<i>nologies, Volume 1 (Long and Short Papers), Min-</i>	Memarrast, Ahmadreza Mosallanezhad, Erfan	717
663	<i>neapolis, Minnesota.</i> Association for Computational	Noury, Shahab Raji, Mohammad Sadegh Rasooli,	718
664	Linguistics.	Sepideh Sadeghi, Erfan Sadeqi Azer, Niloofar Safi	719
665	Zi-Yi Dou, Keyi Yu, and Antonios Anastasopoulos.	Samghabadi, Mahsa Shafaei, Saber Sheybani,	720
666	2019. Investigating meta-learning algorithms for	Ali Tazarv, and Yadollah Yaghoobzadeh. 2021.	721
667	low-resource natural language understanding tasks.	ParsiNLU: A Suite of Language Understanding	722
668	In <i>Proceedings of the 2019 Conference on Empiri-</i>	Challenges for Persian. <i>Transactions of the</i>	723
669	<i>cal Methods in Natural Language Processing and</i>	<i>Association for Computational Linguistics.</i>	724
670	<i>the 9th International Joint Conference on Natural</i>	Tushar Khot, Ashish Sabharwal, and Peter Clark. 2018.	725
671	<i>Language Processing (EMNLP-IJCNLP).</i> Associa-	SciTail: A textual entailment dataset from science	726
672	tion for Computational Linguistics.	question answering. In <i>Thirty-Second AAAI Confer-</i>	727
673	Marie-Pierre Dubuisson and Anil K. Jain. 1994. A	<i>ence on Artificial Intelligence.</i>	728
674	modified hausdorff distance for object matching. In	Gregory Koch, Richard Zemel, Ruslan Salakhutdinov,	729
675	<i>Proceedings of 12th International Conference on</i>	et al. 2015. Siamese neural networks for one-shot	730
676	<i>Pattern Recognition.</i> IEEE.	image recognition. In <i>ICML deep learning work-</i>	731
		<i>shop.</i>	732

733	Bo Liu, Hao Kang, Haoxiang Li, Gang Hua, and Nuno Vasconcelos. 2020. Few-shot open-set recognition using meta-learning . In <i>2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> .	788
734		789
735		790
736		
737	Ilya Loshchilov and Frank Hutter. 2018. Fixing weight decay regularization in adam .	791
738		792
739	Bill MacCartney. 2009. <i>Natural language inference</i> . Stanford University.	793
740		794
741	Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality . In <i>Advances in Neural Information Processing Systems</i> .	795
742		796
743		797
744		798
745		799
746	Ari Morcos, Maithra Raghu, and Samy Bengio. 2018. Insights on representational similarity in neural networks with canonical correlation . In <i>Advances in Neural Information Processing Systems 31</i> . Curran Associates, Inc.	800
747		801
748		802
749		803
750		804
751		805
752	Shikhar Murty, Tatsunori B. Hashimoto, and Christopher Manning. 2021. DReCa: A general task augmentation strategy for few-shot natural language inference . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> .	806
753		807
754		808
755		809
756		
757		
758	Alex Nichol, Joshua Achiam, and John Schulman. 2018. On first-order meta-learning algorithms . <i>arXiv preprint arXiv:1803.02999</i> .	810
759		811
760		812
761	Farhad Nooralahzadeh, Giannis Bekoulis, Johannes Bjerva, and Isabelle Augenstein. 2020. Zero-shot cross-lingual transfer with meta learning . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> . Association for Computational Linguistics.	813
762		814
763		815
764		816
765		817
766		
767	Ankur Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. 2016. A decomposable attention model for natural language inference . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.	818
768		819
769		820
770		821
771		
772		
773	Barun Patra, Joel Ruben Antony Moniz, Sarthak Garg, Matthew R. Gormley, and Graham Neubig. 2019. Bilingual lexicon induction with semi-supervision in non-isometric embedding spaces . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> . Association for Computational Linguistics.	822
774		823
775		824
776		
777		
778		
779		
780	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> . Association for Computational Linguistics.	825
781		826
782		827
783		828
784		829
785		830
786		831
787		
	Matúš Pikuliak, Marián Šimko, and Mária Bielíková. 2021. Cross-lingual learning for text processing: A survey . <i>Expert Systems with Applications</i> .	832
		833
		834
		835
		836
		837
	Kunxun Qi and Jianfeng Du. 2020. Translation-based matching adversarial network for cross-lingual natural language inference . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> .	838
		839
		840
		841
		842
	Kun Qian and Zhou Yu. 2019. Domain adaptive dialog generation via meta learning . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> . Association for Computational Linguistics.	843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

for *Computational Linguistics: Main Volume*. Association for Computational Linguistics.

Zihan Wang, Stephen Mayhew, Dan Roth, et al. 2019. [Cross-lingual ability of multilingual bert: An empirical study](#). *arXiv preprint arXiv:1912.07840*.

Shijie Wu and Mark Dredze. 2019. [Beto, bentz, becas: The surprising cross-lingual effectiveness of BERT](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics.

Mengzhou Xia, Guoqing Zheng, Subhabrata Mukherjee, Milad Shokouhi, Graham Neubig, and Ahmed Hassan Awadallah. 2021. [MetaXL: Meta representation transformation for low-resource cross-lingual learning](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics.

Mingzhang Yin, George Tucker, Mingyuan Zhou, Sergey Levine, and Chelsea Finn. 2019. [Meta-learning without memorization](#). *arXiv preprint arXiv:1912.03820*.

Wenpeng Yin. 2020. [Meta-learning for few-shot natural language processing: A survey](#). *arXiv preprint arXiv:2007.09604*.

A Appendix

A.1 MAML Description

MAML is one of the most popular meta-learning algorithms and it has proven its effectiveness in various fields (e.g., computer vision). MAML is able to find good initialization parameter values and adapt to new tasks quickly. This algorithm can be performed in one episode by following these steps:

- Make a copy of the model with its initial parameters θ .
- Use the training set D_i^{train} to train the model as

$$\hat{\theta} = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(\theta, D_i^{train}) \quad (9)$$

- Apply the model with the updated parameters $\hat{\theta}$ to the validation set D_i^{val} .
- Use the loss on the validation set to update the initial parameters θ

$$\theta = \theta - \beta \nabla_{\theta} \sum_i \mathcal{L}_i(\hat{\theta}, D_i^{val}) \quad (10)$$

A.2 Hyperparameters

Models are implemented using the PyTorch² framework. ParsBERT, mBERT and XLM-R implementations are taken from the HuggingFace library³.

In our experiments, we used the AdamW optimizer (Loshchilov and Hutter, 2018) with learning rate 1e-5 to perform the inner loop of the Reptile algorithm (2), which is known as meta-step. The hyperparameters for the Reptile algorithm are listed in Table 4.

Hyperparameter	Value
epochs	2
number of iterations	20000
sequence length (for NLI)	128
sequence length (for QA)	384
dropout	0.1
optimizer	AdamW
learning rate	1e-5
update steps (m)	3
number of class per episode (way)	2
queue length	4
temperature parameter (τ)	1

Table 4: Hyperparameters for the Reptile algorithm

The hyperparameters for the Prototypical algorithm are also shown in Table 5. Some parameters are calculated based on a grid search, such as Distance Cross-Entropy (DCE) and Cross-Entropy (CE) coefficients, and others are chosen similar to the Reptile algorithm.

Hyperparameter	Value
epochs	2
number of iterations	20000
sequence length (for NLI)	128
dropout	0.1
optimizer	AdamW
learning rate	1e-5
number of class per episode (way)	3
DCE coefficient (λ_1)	1.0
CE coefficient (λ_2)	1.0

Table 5: Hyperparameters for the Prototypical algorithm

The number of iterations parameter varies according to the value of the shot, and is chosen to ensure that all instances in the dataset appear at least once in each epoch.

²<https://pytorch.org/>

³<https://huggingface.co/>

910 **A.3 Prototypical Networks**

911 **A.3.1 Loss Function**

912 As we mentioned in section 2.2, the primary loss
913 function of the Prototypical algorithm is DCE.
914 Since a prototype consists of distribution informa-
915 tion from instances associated with it, the choice of
916 these instances may introduce noise in the learned
917 representation if the neural network is trained only
918 by using the DCE loss. We use CE loss in addition
919 to the DCE loss to make the feature space robust
920 to noise. As a whole, we train the model with a
921 combination of DCE loss and CE loss given by the
922 following equation.

$$923 \text{Loss(overall)} = \lambda_1 \text{Loss}(DCE) + \lambda_2 \text{Loss}(CE) \quad (11)$$

924 **A.3.2 Scenarios**

925 We considered two scenarios for making the
926 episodes. In the first scenario, the model is trained
927 only on auxiliary languages, then fine-tuned using
928 the target language. Therefore, Only auxiliary lan-
929 guages are used to generate support and query sets.
930 An episode of the first scenario is shown in Table 6.

931 In the second scenario, in addition to auxiliary
932 languages, we also used the target language for
933 training. So, the support set is constructed from
934 auxiliary language data and the query set is gener-
935 ated from both auxiliary and Persian language data.
936 Table 7 shows an episode of the second scenario.

937 **A.4 Confusion Matrices**

938 The confusion matrices for the top-performing
939 models (8-shot with four auxiliary languages) is de-
940 picted in Fig. 3 showing the success of this method
941 in improving the accuracy in all classes specially
942 the neutral class.

Example	Category
Support set (or Query set)	
<i>In the midst of this amazing amalgam of cultures is a passion for continuity</i> ⇒ <i>A passion for continuity is not the most important of these cultures</i>	neutral
<i>The river plays a central role in all visits to Paris.</i> ⇒ <i>The river is central to all vacations to Paris</i>	entailment
<i>For the moment, he sought refuge in retreat, and left the room precipitately.</i> ⇒ <i>He stayed put and sat on the floor.</i>	contradiction

Table 6: Example for a 3-way 1-shot episode in the first scenario. In this example we select support set and query set samples from English dataset. As support and query sets are generated similarly, only one set is shown in this table.

Example	Category
Support set	
<i>Recuerda que una vez mencionó que su padre era médico?</i> ⇒ <i>Ella mencionó que su padre era médico hace mucho tiempo</i>	neutral
<i>Dies ist etwas anderes als eine Cantina-Leuchte</i> ⇒ <i>Dies ist sicherlich keine Cantina-Leuchte</i>	entailment
<i>Ensuite, il enfonce un tube respiratoire dans la gorge du patient mort.</i> ⇒ <i>Le patient vit toujours.</i>	contradiction
----- English Translation -----	
You remember her once mentioning that her father was a doctor? ⇒ She mentioned her father being a doctor a long time ago.	neutral
This is something other than a cantina fixture. ⇒ This is certainly not a cantina fixture.	entailment
Next he shoves a breathing tube down the dead patient 's throat . ⇒ The patient is still alive.	contradiction
----- Query set -----	
<i>Une pièce qualifie Frank Lloyd Wright de terrible ingénieur.</i> ⇒ <i>Piece a également déclaré que Wright était un bien meilleur concepteur.</i>	neutral
<i>Sus rápidos oídos captaron el sonido del tren que se acercaba.</i> ⇒ <i>Escuchó que el tren se acercaba rápidamente.</i>	entailment
از قرن دوازدهم به بعد ارقام عربی برای نخستین بار در ایتالیا کاربرد یافت. ⇒ فرانسه اولین کشوری بود که از ارقام عربی استفاده کرد.	contradiction
----- English Translation -----	
A piece calls Frank Lloyd Wright an awful engineer. ⇒ Piece also stated Wright was a much better designer.	neutral
Her quick ears caught the sound of the approaching train. ⇒ She heard the train approaching fast.	entailment
From the twelfth century onwards, Arabic numerals were first used in Italy. ⇒ France was the first country to use Arabic numerals.	contradiction

Table 7: Example for a 3-way 1-shot episode in the second scenario. In this example, the support set samples are selected from French, Spanish, and German datasets, respectively, and the query set samples are selected from French, Spanish, and Persian datasets, respectively.

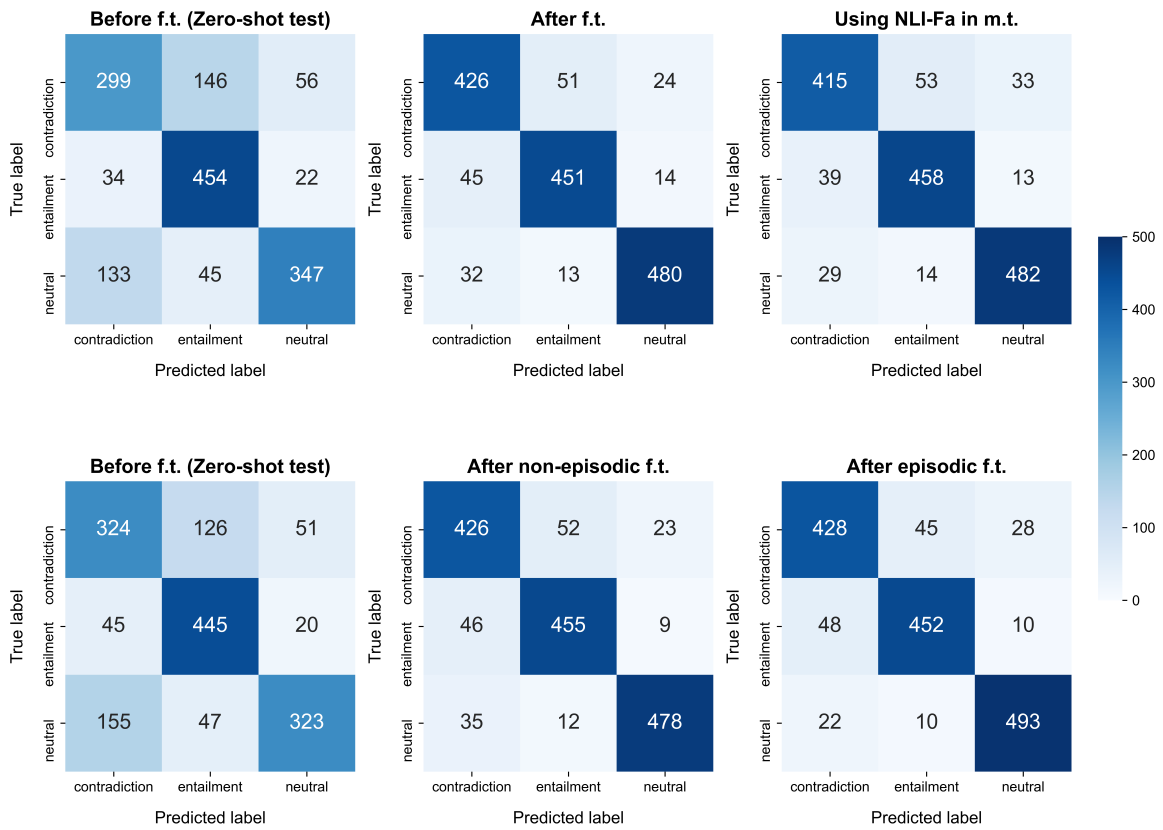


Figure 3: Confusion matrices of the best-obtained model (8-shot with four auxiliary languages) in both meta-learning algorithms on the FarsTail test set. (Top): Reptile algorithm results. (Bottom): Prototypical algorithm results.