# Continual Few-shot Relation Learning via Embedding Space Regularization and Data Augmentation

**Anonymous ACL submission** 

### Abstract

Existing continual relation learning (CRL) methods rely on plenty of labeled training data for learning a new task, which can be hard to acquire in real scenario as getting large and representative labeled data is often expensive and time-consuming. It is therefore necessary for the model to learn novel relational pat-800 terns with very few labeled data while avoiding catastrophic forgetting of previous task knowledge. In this paper, we formulate this challenging yet practical problem as contin-011 ual few-shot relation learning (CFRL). Based 012 on the finding that learning for new emerging few-shot tasks often results in feature dis-014 tributions that are incompatible with previous tasks' learned distributions, we propose a 017 novel method based on embedding space regularization and data augmentation. Our method generalizes to new few-shot tasks and avoids 019 catastrophic forgetting of previous tasks by enforcing extra constraints on the relational embeddings and by adding extra relevant data in a self-supervised manner. With extensive experiments we demonstrate that our method can significantly outperform previous state-of-the-art methods in CFRL task settings.<sup>1</sup>

# 1 Introduction

027

033

040

**Relation Extraction** (RE) aims to detect the relationship between two entities in a sentence, for example, predicting the relation *birthdate* in the sentence "*Kamala Harris* was born in Oakland, California, on *October 20, 1964.*" for the two entities *Kamala Harris* and *October 20, 1964.* It serves as a fundamental step for downstream tasks such as search and question answering (Dong et al., 2015; Yu et al., 2017). Traditionally, RE methods were built by considering a fixed static set of relations (Miwa and Bansal, 2016; Han et al., 2018a). However, similar to entity recognition, RE is also an open-vocabulary problem (Sennrich et al., 2016), where the relation set keeps growing as new relation types emerge with new data. 041

042

043

047

049

052

054

057

059

060

061

062

063

064

065

066

067

069

071

072

073

074

075

076

077

078

079

A potential solution is to formalize RE as Continual Relation Learning or CRL (Wang et al., 2019). In CRL, the model learns relational knowledge through a sequence of tasks, where the relation set changes dynamically from the current task to the next. The model is expected to perform well on both the novel and previous tasks, which is challenging due to the existence of *Catastrophic Forgetting* phenomenon (McCloskey and Cohen, 1989; French, 1999) in continual learning. In this phenomenon, the model forgets previous relational knowledge after learning new relational patterns.

Existing methods to address catastrophic forgetting in CRL can be divided into three categories: (*i*) regularization-based methods, (*ii*) architecturebased methods, and (*iii*) memory-based methods. Recent work shows that memory-based methods which save several key examples from previous tasks to a memory and reuse them when learning new tasks are more effective in NLP (Wang et al., 2019; Sun et al., 2020). Successful memory-based CRL methods include EAEMR (Wang et al., 2019), MLLRE (Obamuyide and Vlachos, 2019), EMAR (Han et al., 2020), and CML (Wu et al., 2021).

Despite their effectiveness, one major limitation of these methods is that they all assume plenty of training data for learning new relations (tasks), which is hard to satisfy in real scenario where continual learning is desirable, as acquiring large labeled datasets for every new relation is expensive and sometimes impractical for quick deployment (*e.g.*, RE from news articles during the onset of an emerging event like Covid-19). In fact, one of the main objectives of continual learning is to quickly adapt to new environments or tasks by exploiting previously acquired knowledge, a hallmark of human intelligence (Lopez-Paz and Ranzato, 2017). If the new tasks are *few-shot*, the existing methods suffer from over-fitting as shown later in our

<sup>&</sup>lt;sup>1</sup>Code and models are available at <redacted>

181

132

experiments (§4). Considering that humans can acquire new knowledge from a handful of examples, it is expected for the models to generalize well on the new tasks with few data. We regard this problem as Continual Few-shot Relation Learning or CFRL (Appendix A.1). Indeed, in relation to CFRL, Zhang et al. (2021), Zhu et al. (2021) and Chen and Lee (2021) recently introduce methods for incremental few-shot learning in Computer Vision.

087

096

100

101

102

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

127

128

129

130

131

Based on the observation that the learning of emerging few-shot tasks may result in distorted feature distributions of new data which is incompatible with previous embedding space (Ren et al., 2020), this work introduces a novel model based on Embedding space Regularization and Data Augmentation (ERDA) for CFRL. In particular, we propose a multi-margin loss and a pairwise margin loss in addition to the traditional cross-entropy loss to impose further relational constraints in the embedding space. We also introduce a novel contrastive loss to learn more effectively from the memory data. Our proposed data augmentation method selects relevant samples from unlabeled text to provide more relational knowledge for the few-shot tasks. The empirical results show that our method can significantly outperform previous state-of-the-art methods. In summary, our main contributions are:

- To the best of our knowledge, we are the first one to consider CFRL. We define the CFRL problem and construct a benchmark for the problem.
- We propose ERDA, a novel method for CFRL based on embedding space regularization and data augmentation.
- With extensive experiments, we demonstrate the effectiveness of our method compared to existing ones and analyse our results thoroughly.

# 2 Related Work

Conventional RE methods include supervised (Zelenko et al., 2002; Liu et al., 2013; Zeng et al., 2014; Miwa and Bansal, 2016), semi-supervised (Chen et al., 2006; Sun et al., 2011; Hu et al., 2020) and distantly supervised methods (Mintz et al., 2009; Yao et al., 2011; Zeng et al., 2015; Han et al., 2018a). These methods rely on a predefined relation set and have limitations in real scenario where novel relations are emerging. There have been some efforts which focus on relation learning without predefined types, including open RE (Shinyama and Sekine, 2006; Etzioni et al., 2008; Cui et al., 2018; Gao et al., 2020) and continual relation learning (Wang et al., 2019; Obamuyide and Vlachos, 2019; Han et al., 2020; Wu et al., 2021).

Continual Learning (CL) aims to learn knowledge from a sequence of tasks. The main problem CL attempts to address is *catastrophic forgetting* (McCloskey and Cohen, 1989), i.e., the model forgets previous knowledge after learning new tasks. Existing methods to alleviate this problem can be divided into three categories. First, regularizationbased methods impose constraints on the update of neural weights important to previous tasks to alleviate catastrophic forgetting (Li and Hoiem, 2017; Kirkpatrick et al., 2017; Zenke et al., 2017; Ritter et al., 2018). Second, architecture-based methods dynamically change model architectures to acquire new information while remembering previous knowledge (Chen et al., 2016; Rusu et al., 2016; Fernando et al., 2017; Mallya et al., 2018). Finally, memory-based methods maintain a memory to save key samples of previous tasks to prevent forgetting (Rebuffi et al., 2017; Lopez-Paz and Ranzato, 2017; Shin et al., 2017; Chaudhry et al., 2019).

**Few-shot Learning** (FSL) aims to solve tasks containing only a few labeled samples, which faces the issue of over-fitting. To address this, existing methods have explored three different directions: (*i*) *data-based* methods use prior knowledge to augment data to the few-shot set (Santoro et al., 2016; Benaim and Wolf, 2018; Gao et al., 2020); (*ii*) *model-based* methods reduce the hypothesis space using prior knowledge (Rezende et al., 2016; Triantafillou et al., 2017; Hu et al., 2018); and (*iii*) *algorithm-based* methods try to find a more suitable strategy to search for the best hypothesis in the whole hypothesis space (Hoffman et al., 2013; Ravi and Larochelle, 2017; Finn et al., 2017).

**Summary.** Existing work in CRL which involves a sequence of tasks containing *sufficient* training data, mainly focuses on alleviating the catastrophic forgetting of previous relational knowledge when the model is trained on new tasks. The work in fewshot learning mostly leverages prior knowledge to address the over-fitting of novel few-shot tasks. In contrast to these lines of work, we aim to solve a more challenging yet more practical problem CFRL where the model needs to learn relational patterns from a sequence of few-shot tasks continually.

# 3 Methodology

In this section, we first formally define the CFRL problem. Then, we present our method for CFRL.

210

211

212

213

214

215

216

218

219

220

221

224

229

231

### **3.1 Problem Definition**

CFRL involves learning from a sequence of tasks  $\mathbb{T} = (\mathcal{T}^1, \dots, \mathcal{T}^n)$ , where every task  $\mathcal{T}^k$  has its own training set  $D_{\text{train}}^k$ , validation set  $D_{\text{valid}}^k$ , and test set  $D_{\text{test}}^k$ . Each dataset D contains several sam-186 ples  $\{(x_i, y_i)\}_{i=1}^{|D|}$ , whose labels  $y_i$  belong to the relation set  $R^k$  of task  $\mathcal{T}^k$ . In contrast to the previously addressed continual relation learning (CRL), CFRL assumes that except for the first task which 190 has enough data for training, the subsequent new 191 tasks are all *few-shot*, meaning that they have only 192 few labeled instances (see Appendix A.1). For 193 example, consider there are three relation learn-194 ing tasks  $\mathcal{T}^1, \mathcal{T}^2$  and  $\mathcal{T}^3$  with their corresponding relation sets  $R^1, R^2$ , and  $R^3$ , each having 10 rela-196 tions. In CFRL, we assume the existing task  $\mathcal{T}^1$  has enough training data (e.g., 100 samples for every 198 relation in  $\mathbb{R}^1$ ), while the new tasks  $\mathcal{T}^2$  and  $\mathcal{T}^3$  are 199 few-shot with only few (e.g., 5) samples for every relation in  $R^2$  and  $R^3$ . Assuming that the relation number of each few-shot task is N and the sample number of every relation is K, we call this setup Nway K-shot continual learning. The problem setup of CFRL is aligned with the real scenario, where we generally have sufficient data for an existing task, but only few labeled data as new tasks emerge.

> The model in CFRL is expected to first learn  $\mathcal{T}^1$ well, which has sufficient training data to obtain good ability to extract the relation information in the sentence. Then at time step k, the model will be trained on the training set  $D_{\text{train}}^k$  of few-shot task  $\mathcal{T}^k$ . After learning  $\mathcal{T}^k$ , the model is expected to perform well on both  $\mathcal{T}^k$  and the previous k-1tasks, as the model will be evaluated on  $\hat{D}_{\text{test}}^k = \bigcup_{i=1}^k D_{\text{test}}^i$  consisting of all known relations after learning  $\mathcal{T}^k$ , *i.e.*,  $\hat{R}^k = \bigcup_{i=1}^k R^i$ . This requires the model to overcome the *catastrophic forgetting* of previous knowledge and to learn new knowledge well with very few labeled data.

To overcome the catastrophic forgetting problem, a memory  $\mathcal{M} = \{\mathcal{M}^1, \mathcal{M}^2, ...\}$ , which stores some key samples of previous tasks is maintained during the learning. When the model is learning  $\mathcal{T}^k$ , it has access to the data saved in memory  $\mathcal{M}^1, ..., \mathcal{M}^{k-1}$ . As there is no limit on the number of tasks, the size of memory  $\mathcal{M}^k$  is constrained to be small. Therefore, the model has to select only key samples from the training set  $D_{\text{train}}^k$  to save them in  $\mathcal{M}^k$ . In our CFRL setting, only one sample per relation is allowed to be saved in the memory.



Figure 1: Our framework for CFRL. The Data Augmentation component is used only for few-shot tasks (k > 1).

232

233

234

237

238

240

241

242

243

245

246

247

249

250

251

252

253

254

255

257

258

259

260

261

262

265

267

269

270

271

# 3.2 Overall Framework

Our framework for CFRL is shown in Fig. 1 and Alg. 1 describes the overall training process (see Appendix A.2 for a block diagram). At time step k, given the training data  $D_{\text{train}}^k$  for the task  $\mathcal{T}^k$ , depending on whether the task is a few-shot or not, the process has four or three working modules, respectively. The general learning process (§3.3) has three steps that apply to all tasks. If the task is a few-shot task (k > 1), we apply an additional step to create an augmented training set  $\tilde{D}_{\text{train}}^k$ . For the initial task (k = 1), we have  $\tilde{D}_{\text{train}}^k = D_{\text{train}}^k$ . For any task  $\mathcal{T}^k$ , we use a siamese model to en-

For any task  $\mathcal{T}^k$ , we use a siamese model to encode every new relation  $r_i \in \mathbb{R}^k$  into  $\mathbf{r}_i \in \mathbb{R}^d$ as well as the sentences, and train the model on  $\widetilde{D}_{\text{train}}^k$  to acquire relation information of the new data (§3.3.2). To overcome forgetting, we select the most informative sample for each relation  $r_i \in \mathbb{R}^k$ from  $D_{\text{train}}^k$  and update the memory  $\widehat{\mathcal{M}}^k$  (§3.3.3). Finally, we combine  $\widetilde{D}_{\text{train}}^k$  and  $\widehat{\mathcal{M}}^k$  as the training data for learning new relational patterns and remembering previous knowledge (§3.3.4). We also simultaneously update the representation of all relations in  $\widehat{\mathbb{R}}^k$ , which involves making a forward pass through the current model. The learning and updating are done iteratively for convergence.

For data augmentation in the few-shot tasks (§3.4), we select reliable samples with high relational similarity score from an unlabelled corpus C using a fine-tuned BERT (Devlin et al., 2019), which serves as the relational similarity model  $S_{\pi}$ . In the interests of coherence, we first present the general learning method followed by the augmentation process for few-shot learning.

# 3.3 General Learning Process

We first introduce the encoder network as it is the basic component of the whole framework.

### 3.3.1 The Encoder Network

The siamese encoder  $(f_{\theta})$  aims at extracting generic and relation related features from the input. The **Require:** the training set  $D_{\text{train}}^k$  and the relation set  $R^k$  of the current task  $\mathcal{T}^k$ , the current memory  $\hat{\mathcal{M}}^{k-1}$  and the known relation set  $\hat{R}^{k-1}$ , the model  $\theta$ , the similarity model  $S_{\pi}$ , and the unlabeled text corpus C.

1: if 
$$k == 1$$
 then  $\triangleright$  initial task  
2:  $\widetilde{D}_{\text{train}}^k = D_{\text{train}}^k$ 

- 3: else ⊳ few-shot task SELECT similar samples from C using  $S_{\pi}$  for every 4:
- sample in  $D_{\text{train}}^k$  and store them in A
- $\widetilde{D}_{\text{train}}^k = A \cup D_{\text{train}}^k$ 5:
- 6: end if

272

273

274

275

276

278

279

281

283

287

289

290

291

292

296

297

- 7: INITIALIZE  $\mathbf{r}_i$  for every relation  $r_i \in R^k$
- 8: for  $i = 1, ..., iter_1$  do
- UPDATE  $\theta$  with  $\mathcal{L}_{\text{new}}$  on  $\widetilde{D}_{\text{train}}^k$ ▷ Train on new task 9: 10: end for
- 11: SELECT key samples from  $D_{\text{train}}^k$  for every relation  $r_i \in$
- $R^{k} \text{ to save in } \mathcal{M}^{k}$ 12:  $\hat{R}^{k} = \hat{R}^{k-1} \cup R^{k}$ 13:  $\hat{\mathcal{M}}^{k} = \hat{\mathcal{M}}^{k-1} \cup \mathcal{M}^{k}$ ▷ Update memory 14:  $\widetilde{H}^k = \widetilde{D}^k_{\text{train}} \cup \mathcal{\hat{M}}^k$ ▷ Combine two data sources 15: for  $i = 1, ..., iter_2$  do UPDATE heta with  $\mathcal{L}_{\text{mem}}$  on  $\widetilde{H}^k$ 16: UPDATE  $\mathbf{r}_i$  for every relation  $r_i \in \hat{R}^k$ 17: 18: end for
- input can be a labeled sentence or the name of a relation. We adopt two kinds of encoders:

**Bi-LSTM** To have a fair comparison with previous work, we use the same architecture as Han et al. (2020). It takes GloVe embeddings (Pennington et al., 2014) of the words in a given input and produces a vector representation through a Bi-LSTM (Hochreiter and Schmidhuber, 1997).

**BERT** We adopt BERT<sub>base</sub> which has 12 layers and 110M parameters. As the new tasks are fewshot, we only fine-tune the 12-th encoding layer and the extra linear layer. We include special tokens around the entities ('#' for the head entity and '@' for the tail entity) in a given labeled sentence to improve the encoder's understanding of relation information. We use the [CLS] token features as the representation of the input sequence.

### 3.3.2 Learning with New Data

At time step k, to have a good understanding of the new relations, we fine-tune the model on the expanded dataset  $\widetilde{D}_{\text{train}}^k$ . The model  $f_{\theta}$  first encodes the name of each new relation  $r_i \in \mathbb{R}^k$  into its representation  $\mathbf{r}_i \in \mathbb{R}^d$  by making a forward pass. Then, we optimize the parameters  $(\theta)$  by minimizing a loss  $\mathcal{L}_{new}$  that consists of a cross entropy loss, a multi-margin loss and a pairwise margin loss.

The **cross entropy** loss  $\mathcal{L}_{ce}$  is used for relation

classification as follows.

$$-\sum_{(x_i,y_i)\in \tilde{D}_{\text{train}}^k} \sum_{j=1}^{|R^k|} \delta_{y_i,r_j} \times \log \frac{\exp(g(f_\theta(x_i),\mathbf{r}_j))}{\sum_{l=1}^{|\hat{R}^k|} \exp(g(f_\theta(x_i),\mathbf{r}_l))}$$
(1)

where  $\hat{R}^k$  is the set of all known relations at step k, q(,) is a function used to measure similarity between two vectors (e.g., cosine similarity or L2 distance), and  $\delta_{a,b}$  is the Kronecker delta function–  $\delta_{a,b} = 1$  if a equals b, otherwise  $\delta_{a,b} = 0$ .

In inference, we choose the relation label that has the highest similarity with the input sentence (Eq. 8). To ensure that an example has the highest similarity with the true relation, we additionally design two margin-based losses, which increase the score between an example and the true label while decreasing the scores for the wrong labels. The first one is a multi-margin loss defined as:

$$\mathcal{L}_{mm} = \sum_{(x_i, y_i) \in \widetilde{D}_{train}^k} \sum_{j=1, j \neq t_i}^{|\widehat{R}^k|} \max\left(0, \\ m_1 - g(f_{\theta}(x_i), \mathbf{r}_{t_i}) + g(f_{\theta}(x_i), \mathbf{r}_j)\right)$$
(2)

where  $t_i$  is the correct relation index in  $\hat{R}^k$  satisfying  $r_{t_i} = y_i$  and  $m_1$  is a margin value. The  $\mathcal{L}_{mm}$  loss attempts to ensure intra-class compactness while increasing inter-class distances. The second one is a **pairwise margin** loss  $\mathcal{L}_{pm}$ :

$$\sum_{i_i, y_i) \in \widetilde{D}_{\text{train}}^k} \max\left( 0, m_2 - g(f_\theta(x_i), \mathbf{r}_{t_i}) + g(f_\theta(x_i), \mathbf{r}_{s_i}) \right)$$
(3)

where  $m_2$  is the margin for  $\mathcal{L}_{pm}$  and  $s_i =$  $\arg \max_{s} g(f_{\theta}(x_i), \mathbf{r}_s)$  s.t.  $s \neq t_i$ , the closest wrong label. The  $\mathcal{L}_{pm}$  loss penalizes the cases where the similarity score of the closest wrong label is higher than the score of the correct label (Yang et al., 2018). Both  $\mathcal{L}_{mm}$  and  $\mathcal{L}_{pm}$  improve the discriminative ability of the model (§4.4).

The **total loss** for learning on  $\mathcal{T}^k$  is defined as:

$$\mathcal{L}_{new} = \lambda_{ce} \mathcal{L}_{ce} + \lambda_{mm} \mathcal{L}_{mm} + \lambda_{pm} \mathcal{L}_{pm} \qquad (4)$$

where  $\lambda_{ce}$ ,  $\lambda_{mm}$  and  $\lambda_{pm}$  are the relative weights of the component losses, respectively.

### **3.3.3** Selecting Samples for Memory

After training the model  $f_{\theta}$  with Eq. (4), we use it to select one sample per new relation. Specifically, for every new relation  $r_j \in R^k$ , we obtain the centroid feature  $c_i$  by averaging the embeddings of all samples labeled as  $r_i$  in  $D_{\text{train}}^k$  as follows.

$$\mathbf{c}_j = \frac{1}{|D_{r_j}^k|} \sum_{(x_i, y_i) \in D_{r_j}^k} f_\theta(x_i) \tag{5}$$

(x

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

327

328

330

331

332

333

334

335

336

337

339 where  $D_{r_j}^k = \{(x_i, y_i) | (x_i, y_i) \in D_{\text{train}}^k, y_i = r_j\}.$ 340 Then we select the instance closest to  $\mathbf{c}_j$  from  $D_{r_j}^k$ 341 as the most informative sample and save it in mem-342 ory  $\mathcal{M}^k$ . Note that the selection is done from  $D_{\text{train}}^k$ , 343 not from the expanded set  $\widetilde{D}_{\text{train}}^k$ .

345

347

354

356

361

365

367

371

374

376

**3.3.4** Alleviating Forgetting through Memory As the learning of new relational patterns may cause catastrophic forgetting of previous knowledge (see baselines in §4), our model needs to learn from the memory data to alleviate forgetting. We combine the expanded set  $\widetilde{D}_{train}^k$  and the whole memory data  $\widehat{\mathcal{M}}^k = \bigcup_{j=1}^k \mathcal{M}^j$  into  $\widetilde{H}^k$  to allow the model to learn new relational knowledge and consolidate previous knowledge. However, the memory data is limited containing only one sample per relation. To learn effectively from such limited data, we design a novel method to generate a *hard negative sample* set  $P_i$  for every sample in  $\widehat{\mathcal{M}}^k$ .

The negative samples are generated on the fly. After sampling a mini-batch  $B_t$  from  $\tilde{H}^k$ , we consider all memory data in  $B_t$  as  $M_{B_t}$ . For every sample  $(\hat{x}_i, \hat{y}_i)$  in  $M_{B_t}$ , we replace its head entity  $e_i^h$  or tail entity  $e_i^t$  with the corresponding entity of a randomly selected sample in the same batch  $B_t$  to get the hard negative sample set  $P_i = \{(\hat{x}_j^{P_i}, \hat{y}_i)\}_{j=1}^{|P_i|}$ . Then  $(\hat{x}_i, \hat{y}_i)$  and  $P_i$  are used to calculate a margin-based **contrastive loss**  $\mathcal{L}_{con}$  as follows.

$$\mathcal{L}_{con} = \sum_{(\hat{x}_{i}, \hat{y}_{i}) \in M_{B_{t}}} \max\left(0, m_{3} - g(f_{\theta}(\hat{x}_{i}), \mathbf{r}_{\hat{t}_{i}}) + \sum_{(\hat{x}_{j}^{P_{i}}, \hat{y}_{i}) \in P_{i}} g(f_{\theta}(\hat{x}_{j}^{P_{i}}), \mathbf{r}_{\hat{t}_{i}})\right)$$
(6)

where  $\hat{t}_i$  is the relation index satisfying  $r_{\hat{t}_i} = \hat{y}_i$ and  $m_3$  is the margin value for  $\mathcal{L}_{con}$ . This loss forces the model to distinguish the valid relations from the hard negatives so that the model learns more precise and fine-grained relational knowledge. In addition, we also use the three losses  $\mathcal{L}_{ce}$  and  $\mathcal{L}_{mm}$  and  $\mathcal{L}_{pm}$  defined in §3.3.2 to update  $\theta$  on  $B_t$ . The total loss on the memory data is:

$$\mathcal{L}_{mem} = \lambda_{ce} \mathcal{L}_{ce} + \lambda_{mm} \mathcal{L}_{mm} + \lambda_{pm} \mathcal{L}_{pm} + \lambda_{con} \mathcal{L}_{con}$$
(7)

where  $\lambda_{ce}$ ,  $\lambda_{mm}$ ,  $\lambda_{pm}$  and  $\lambda_{con}$  are the relative weights of the corresponding losses.

378Updating Relation EmbeddingsAfter training379the model on  $\tilde{H}^k$  for few steps, we use the mem-380ory  $\hat{\mathcal{M}}^k$  to update the relation embedding  $\mathbf{r}_i$  of all381known relations. For a relation  $r_i \in \hat{R}^k$ , we aver-382age the embeddings (obtained by making a forward

pass through  $f_{\theta}$ ) of the relation name and memory data to obtain its updated representation  $\mathbf{r}_i$ . The training of  $\theta$  and updating of  $\mathbf{r}_i$  is done iteratively to grasp new relational patterns while alleviating the catastrophic forgetting of previous knowledge.

383

384

386

388

391

392

393

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

# 3.3.5 Inference

For a given input  $x_i$  in  $\hat{D}_{\text{test}}^k$ , we calculate the similarity between  $x_i$  and all known relations, and pick the one with the highest similarity score:

$$y_i^* = \operatorname*{arg\,max}_{r \in \hat{R}^k} g(f_\theta(x_i), \mathbf{r}) \tag{8}$$

### **3.4 Data Augmentation for Few-shot Tasks**

For each few-shot task  $\mathcal{T}^k$ , we aim to get more data by selecting reliable samples from an unlabeled corpus  $\mathcal{C}$  with tagged entities before the general learning process (§3.3) begins. We achieve this using a relational similarity model  $S_{\pi}$  and sentences from WikiPedia as  $\mathcal{C}$ . The model  $S_{\pi}$  (described later) takes a sentence as input and produces a normalized vector representation. The cosine similarity between two vectors is used to measure the relational similarity between the two corresponding sentences. A higher similarity means the two sentences are more likely to have the same relation label. We propose two novel selection methods, which are complementary to each other.

(a) Augmentation via Entity Matching For each instance  $(x_i, y_i)$  in  $D_{\text{train}}^k$ , we extract its entity pair  $(e_i^h, e_i^t)$  with  $e_i^h$  being the head entity and  $e_i^t$  being the tail entity. As sentences with the same entity pair are more likely to express the same relation, we first collect a candidate set  $\mathcal{Q} = \{\tilde{x}_j\}_{j=1}^{|\mathcal{Q}|}$  from  $\mathcal{C}$ , where  $\tilde{x}_j$  shares the same entity pair  $(e_i^h, e_i^t)$  with  $x_i$ . If  $\mathcal{Q}$  is a non-empty set, we pair all  $\tilde{x}_j$  in  $\mathcal{Q}$ with  $x_i$ , and denote each pair as  $\langle \tilde{x}_j, x_i \rangle$ . Then we use  $S_{\pi}$  to obtain a similarity score  $s_j$  for  $\langle \tilde{x}_j, x_i \rangle$ . After getting scores for all pairs, we pick the instances  $\tilde{x}_j$  with similarity score  $s_j$  higher than a predefined threshold  $\alpha$  as new samples and label them with relation  $y_i$ . The selected instances are then augmented to  $D_{\text{train}}^k$  as additional data.

(b) Augmentation via Similarity Search The hard entity matching could be too restrictive at times. For example, even though the sentences *"Harry Potter* is written by *Joanne Rowling"* and *"Charles Dickens* is the author of *A Tale of Two Cities"* share the same relation *author*, hard matching fails to find any relevance. Therefore, in cases

when entity matching returns an empty Q, we re-430 sort to similarity search using Faiss (Johnson et al., 431 2017). Given a query vector  $q_i$ , it can efficiently 432 search for vectors  $\{\mathbf{v}_j\}_{j=1}^K$  with the top-K highest 433 similarity scores in a large vector set  $\mathcal{V}$ . In our case, 434  $\mathbf{q}_i$  is the representation of  $x_i$  and  $\mathcal{V}$  contains the 435 representations of the sentences in C. We use  $S_{\pi}$ 436 to obtain these representations; the difference is 437 that  $\mathcal{V}$  is pre-computed while  $\mathbf{q}_i$  is obtained dur-438 ing training. We labeled the top-K most similar 439 instances with  $y_i$  and augment them to  $D_{\text{train}}^k$ . 440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

**Similarity Model** To train  $S_{\pi}$ , inspired by Soares et al. (2019), we adopt a *contrastive learning* method to fine-tune a BERT<sub>base</sub> model on C, whose sentences are already tagged with entities. Based on the observation that sentences with the same entity pair are more likely to encode the same relation, we use sentence pairs containing the same entities in C as positive samples. For negatives, instead of using all sentence pairs containing different entities, we select pairs sharing only one entity as **hard negatives** (*i.e.*, pair  $(x_i, x_j)$  where  $e_i^h = e_j^h$  and  $e_i^t \neq e_j^t$  or  $e_i^t = e_j^t$  and  $e_i^h \neq e_j^h$ ). We randomly sample the same number of negative samples as the positive ones to balance the training.

For an input pair  $(x_i, x_j)$ , we compute the similarity score based on the following formula.

$$\sigma(x_i, x_j) = \frac{1}{1 + \exp(-\mathcal{S}_{\pi}(x_i)^T \mathcal{S}_{\pi}(x_j))} \quad (9)$$

where  $S_{\pi}(x)$  is the normalized representation of xobtained from the final layer of BERT. Then we optimize the parameters  $\pi$  of  $S_{\pi}$  by minimizing a binary cross entropy loss  $\mathcal{L}_{\text{pretrain}}$  as follows.

$$-\sum_{(x_i,x_j)\in\mathcal{C}_p}\log\sigma(x_i,x_j) - \sum_{(x'_i,x'_j)\in\mathcal{C}_n}\log(1-\sigma(x'_i,x'_j)) \quad (10)$$

where  $C_p$  is a positive batch and  $C_n$  is a negative batch. This objective tries to ensure that sentence pairs with the same entity pairs have higher cosine similarity than those with different entities.

# 4 Experiment

We define the benchmark and evaluation metric for CFRL before presenting our experimental results.

# 4.1 Benchmark and Evaluation Metric

**Benchmark** As the benchmark for CFRL needs to have sufficient relations as well as data and be suitable for few-shot learning, we create the CFRL benchmark based on **FewRel** (Han et al., 2018b). FewRel is a large-scale dataset for few-shot RE, which contains 80 relations with hundreds of samples per relation. We randomly split the 80 relations into 8 tasks, where each task contains 10 relations (10-way). To have enough data for the first task  $T^1$ , we sample 100 samples per relation. All the subsequent tasks  $T^2$ , ...,  $T^8$  are few-shot; for each relation, we conduct 2-shot, 5-shot and 10-shot experiments to verify the effectiveness of our method. 475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

In addition, to demonstrate the generalizability of our method, we also create a CFRL benchmark based on the **TACRED** dataset (Zhang et al., 2017) which contains only 42 relations. We filter out the special relation "n/a" (not available) and split the remaining 41 relations into 8 tasks. Except for the first task that contains 6 relations, all other tasks have 5 relations (5-way). Similar to FewRel, we randomly sample 100 examples per relation in  $\mathcal{T}^1$ and conduct 5-shot and 10-shot experiments.

**Metric** At time step k, we evaluate the model performance through relation classification accuracy on the test sets  $\hat{D}_{\text{test}}^k = \bigcup_{i=1}^k D_{\text{test}}^i$  of all seen tasks  $\{\mathcal{T}^i\}_{i=1}^k$ . This metric reflects whether the model can alleviate catastrophic forgetting while acquiring novel knowledge well with very few data. Since the model performance might be influenced by task sequences and few-shot training samples, we run every experiment 6 times with different random seeds and report the average accuracy.

# 4.2 Model Settings & Baselines

We follow the settings in (Han et al., 2020) for the Bi-LSTM encoder to have a fair comparison. For data augmentation, we set the threshold  $\alpha = 0.65$ and the number of samples selected by Faiss (K) as 1. We adopt 0.2, 0.2 and 0.01 for the three margin values  $m_1, m_2$  and  $m_3$ , respectively. The loss weights  $\lambda_{ce}, \lambda_{mm}, \lambda_{pm}$  and  $\lambda_{con}$  are set to 1.0, 1.0, 1.0 and 0.01, respectively. In Alg. 1, we set 1 for *iter*<sub>1</sub> and 2 for *iter*<sub>2</sub>. Hyperparameter search was done on the validation sets (Appendix A.3). We compare our approach with the following baselines:

- SeqRun fine-tunes the model only on the training data of the new tasks without using any memory data. It may face serious catastrophic forgetting and serves as a lower bound.
- Joint Training stores all previous samples in the memory and trains the model on all data for each new task. It serves as an **upper bound** in **CRL**.
- EMR (Wang et al., 2019) maintains a memory for storing selected samples from previous tasks.

Method	Task index								
	1	2	3	4	5	6	7	8	
SeqRun Joint Training	92.78 <b>92.78</b>	52.11 76.29	30.08 69.39	24.33 64.75	19.83 60.45	16.90 <b>57.64</b>	14.36 52.80	12.34 50.03	
EMR EMAR IDLVQ-C	92.78 85.20 92.23	69.14 62.02 69.15	56.24 52.45 57.42	50.03 48.95 51.66	46.50 46.77 49.31	43.21 44.33 46.24	39.88 40.75 42.25	37.51 39.04 40.56	
ERDA ERDA(BERT)	92.57 94.22	<b>79.17</b> 87.72	<b>70.43</b> 82.66	<b>65.01</b> 78.29	<b>61.06</b> 73.99	<b>57.54</b> 69.45	<b>54.88</b> 68.08	<b>53.23</b> 65.30	

Table 1: Accuracy (%) of different methods at every time step on **FewRel** benchmark for 10-way 5-shot CFRL. ERDA represents our method with a Bi-LSTM encoder and ERDA(BERT) represents our method with a BERT encoder.



Figure 2: Comparison of the results at each time step on the **FewRel** benchmark for 10-way 2-shot and 10-shot settings.

When training on a novel task, EMR combines the new training data and memory data.

- EMAR (Han et al., 2020) is the state-of-the-art on CRL, which adopts memory activation and reconsolidation to alleviate catastrophic forgetting.
- **IDLVQ-C** (Chen and Lee, 2021) introduces quantized reference vectors to represent previous knowledge and mitigates catastrophic forgetting by imposing constraints on the quantized vectors and embedded space. It was originally proposed for image classification with state-of-the-art results in incremental few-shot learning.

#### 4.3 Main Results

**FewRel Benchmark** We report our results on 10way 5-shot in Table 1, while Fig. 2 shows the results on the 10-way 2-shot and 10-way 10-shot settings. From the results, we can observe that:

• Our proposed ERDA outperforms previous baselines in all CFRL settings, which demonstrates the superiority of our method. Simply fine-tuning the model with new few-shot examples leads to rapid drops in accuracy due to severe over-fitting and catastrophic forgetting. Although EMR and EMAR adopt a memory module to alleviate forgetting, their performance still decreases quickly as they require plenty of training data for learning a new task. Compared with EMR and EMAR, IDLVQ-C is slightly better as it introduces quantized vectors



Figure 3: t-SNE visualization of IDLVQ-C and ERDA at two stages. Colors represent different relation classes with numbers being the relation indices. The initial embeddings of four base classes after learning the first task are shown in the upper row. As the data for the first task is sufficient, both methods can obtain separable embedding space. The lower row shows the embeddings of four base classes and two novel classes (Id 5 and 9) after learning a new few-shot task. Compared with IDLVQ-C, ERDA shows better intra-class compactness (circled regions) and larger inter-class distances (see the distances between 5 and 9, and 9 and 65).

553

554

555

556

557

559

560

561

562

563

564

565

566

567

568

570

571

572

573

574

575

576

578

579

580

582

that can better represent the embedding space of few-shot tasks. However, IDLVQ-C does not necessarily push the samples from different relations to be far apart in the embedding space and the updating method for the reference vectors may not be optimal. ERDA outperforms IDLVQ-C by a large margin through embedding space regularization and self-supervised data augmentation. To verify this, we show the embedding space of IDLVQ-C and ERDA using t-SNE (Van der Maaten and Hinton, 2008). We randomly choose four classes from the first task of FewRel and two classes from the new task, and visualize the test data of these classes in Fig. 3. As can be seen, the embedding space obtained by ERDA shows better intra-class compactness and larger inter-class distances.

• Unlike CRL, joint training does not always serve as an upper bound in CFRL due to the extremely imbalanced data distribution. Benefiting from the ability to learn feature distribution with very few data, both ERDA and IDLVQ-C perform better than joint training in the 2-shot setting. However, as the number of few-shot samples increases, the performance of IDLVQ-C falls far behind joint training, while ERDA still performs better. In the 5-shot setting, ERDA could achieve better results than joint training which verifies the effectiveness of self-supervised data augmentation (more on this in §4.4). Although ERDA performs worse than joint training in the 10-shot setting, its results are

525

527

53

541

546

548

549



Figure 4: Comparison results at every time step on TA-CRED benchmark for 5-way 5-shot and 10-shot settings.

still much better than other baselines.

583

584

585

586

588

590

595

596

599

601

603

604

605

609

610

611

• After learning all few-shot tasks, ERDA outperforms IDLVQ-C by **9.69**%, **12.67**% and **11.49**% in the 2-shot, 5-shot and 10-shot settings, respectively. Moreover, the relative gain of ERDA keeps growing with the increasing number of new few-shot tasks. This demonstrates the ability of our method in handling a longer sequence of CFRL tasks. In addition, compared with ERDA, ERDA(BERT) achieves much better performance. At the last time step, ERDA(BERT) outperforms ERDA by **12.07**%, which demonstrates the power of using pre-trained language models for CFRL.

**TACRED Benchmark** Fig. 4 shows the *5-way 5-shot* and *5-way 10-shot* results on TACRED. We can see that here also ERDA outperforms all other methods by a large margin which verifies the strong generalization ability of our proposed method.

### 4.4 Ablation Study

We conduct several ablations to analyze the contribution of different components of ERDA on the FewRel 10-way 5-shot setting. In particular, we investigate six other variants of ERDA by removing one component at a time: (a) the multi-margin loss  $\mathcal{L}_{mm}$ , (b) the pairwise margin loss  $\mathcal{L}_{pm}$ , (c) the margin-based contrastive loss  $\mathcal{L}_{con}$ , (d) the whole 2-stage data augmentation module, (e) the entity matching method of augmentation, and (f) the similarity search method of augmentation.

From the results in Table 2, we can observe that 612 all components improve the performance of our 613 model. Specifically,  $\mathcal{L}_{mm}$  yields about 1.51% performance boost as it brings samples of the same 615 relation closer to each other while enforcing larger 616 distances among different relation distributions. 617 The  $\mathcal{L}_{pm}$  improves the accuracy by **3.18**%, which demonstrates the effect of contrasting with the near-619 est wrong label. The adoption of  $\mathcal{L}_{con}$  leads to 620 1.28% improvement, which shows that generating 621 hard negative samples for memory data can help to better remember previous relational knowledge.

Method	Task index									
	1	2	3	4	5	6	7	8		
ERDA	92.57	79.17	70.43	65.01	61.06	57.54	54.88	53.23		
w.o. $\mathcal{L}_{mm}$	91.67	78.38	70.21	63.77	60.23	56.32	53.45	51.72		
w.o. $\mathcal{L}_{pm}$	91.37	75.80	67.11	61.13	57.14	54.04	51.59	50.05		
w.o. $\mathcal{L}_{con}$	91.63	79.05	69.28	63.86	59.66	56.68	54.12	51.95		
w.o. DA	92.57	77.84	69.76	63.74	58.31	56.12	53.21	51.51		
w.o. EM	92.57	78.33	70.17	64.18	59.63	57.14	54.18	52.39		
w.o. SS	92.57	78.56	69.94	63.98	59.85	56.92	53.75	52.27		

Table 2: Ablations on FewRel benchmark (10-way 5-shot).

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

The data augmentation module improves the performance by **1.72**% as it can extract informative samples from unlabeled text which provide more relational knowledge for few-shot tasks. The results of variants without entity matching or similarity search verify that the two data augmentation methods are generally complementary to each other.

One could argue that the data augmentation module increases the complexity of ERDA compared to other models. However, astute readers can find that even without data augmentation, ERDA outperforms IDLVQ-C significantly for all tasks (compare 'ERDA *w.o.* DA' with the baselines in Table 1).

**ERDA's Performance under CRL** Although ERDA is designed for CFRL, we also evaluate the embedding space regularization ('ERDA *w.o.* DA') on the CRL setting. We sample 100 examples per relation for every task in FewRel and compare our method with the state-of-the-art method EMAR. The results are shown in Appendix A.4. We can see that ERDA outperforms EMAR in all tasks by **1.01 - 5.90**% proving that the embedding regularization can be a general method for CRL.

# 5 Conclusion

We have introduced continual few-shot relation learning (CFRL), a challenging yet practical problem where the model needs to learn new relational knowledge with very few labeled data continually. We have proposed a novel method, named ERDA, to alleviate the over-fitting and catastrophic forgetting problems which are the core issues in CFRL. ERDA imposes relational constraints in the embedding space with innovative losses and adds extra informative data for few-shot tasks in a selfsupervised manner to better grasp novel relational patterns and remember previous knowledge. Extensive experimental results and analysis show that ERDA significantly outperforms previous methods in all CFRL settings investigated in this work. In the future, we would like to investigate ways to combine meta-learning with CFRL.

### References

665

666

673

674

678

679

681

686

702

709

710

711

714

715

716

718

719

- Sagie Benaim and Lior Wolf. 2018. One-shot unsupervised cross domain translation. In Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pages 2108–2118.
- Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. 2019. Efficient lifelong learning with A-GEM. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net.
- Jinxiu Chen, Donghong Ji, Chew Lim Tan, and Zhengyu Niu. 2006. Relation extraction using label propagation based semi-supervised learning. In Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics, pages 129–136, Sydney, Australia. Association for Computational Linguistics.
- Kuilin Chen and Chi-Guhn Lee. 2021. Incremental few-shot learning via vector quantization in deep embedded space. In *International Conference on Learning Representations*.
- Tianqi Chen, Ian J. Goodfellow, and Jonathon Shlens. 2016. Net2net: Accelerating learning via knowledge transfer. In 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings.
- Lei Cui, Furu Wei, and Ming Zhou. 2018. Neural open information extraction. In *Proceedings of the* 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 407–413, Melbourne, Australia. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Li Dong, Furu Wei, Ming Zhou, and Ke Xu. 2015. Question answering over Freebase with multicolumn convolutional neural networks. In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 260–269, Beijing, China. Association for Computational Linguistics.

Oren Etzioni, Michele Banko, Stephen Soderland, and Daniel S Weld. 2008. Open information extraction from the web. *Communications of the ACM*, 51(12):68–74. 720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

749

750

751

752

753

754

755

756

757

758

759

760

762

763

764

765

766

768

769

770

771

772

774

776

- Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A Rusu, Alexander Pritzel, and Daan Wierstra. 2017. Pathnet: Evolution channels gradient descent in super neural networks. *arXiv preprint arXiv:1701.08734*.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017, volume 70 of Proceedings of Machine Learning Research, pages 1126–1135. PMLR.
- Robert M French. 1999. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4):128–135.
- Tianyu Gao, Xu Han, Ruobing Xie, Zhiyuan Liu, Fen Lin, Leyu Lin, and Maosong Sun. 2020. Neural snowball for few-shot relation learning. In *The Thirty-Fourth AAAI Conference on Artificial Intelli*gence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020, pages 7772– 7779. AAAI Press.
- Xu Han, Yi Dai, Tianyu Gao, Yankai Lin, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. 2020. Continual relation learning via episodic memory activation and reconsolidation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6429–6440, Online. Association for Computational Linguistics.
- Xu Han, Pengfei Yu, Zhiyuan Liu, Maosong Sun, and Peng Li. 2018a. Hierarchical relation extraction with coarse-to-fine grained attention. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2236–2245, Brussels, Belgium. Association for Computational Linguistics.
- Xu Han, Hao Zhu, Pengfei Yu, Ziyun Wang, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2018b. FewRel: A large-scale supervised few-shot relation classification dataset with state-of-the-art evaluation. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 4803– 4809, Brussels, Belgium. Association for Computational Linguistics.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Judy Hoffman, Eric Tzeng, Jeff Donahue, Yangqing Jia, Kate Saenko, and Trevor Darrell. 2013. Oneshot adaptation of supervised deep convolutional models. *arXiv preprint arXiv:1312.6204*.

887

888

833

834

835

Xuming Hu, Fukun Ma, Chenyao Liu, Chenwei Zhang, Lijie Wen, and Philip S Yu. 2020. Semisupervised relation extraction via incremental meta self-training. *Update*, 9:8.

778

790

792

793

794

799

803

804

807

810

811

812

813

814

815

816

817

818

819

821

822

823 824

827

830

831

- Zikun Hu, Xiang Li, Cunchao Tu, Zhiyuan Liu, and Maosong Sun. 2018. Few-shot charge prediction with discriminative legal attributes. In *Proceedings* of the 27th International Conference on Computational Linguistics, pages 487–498, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. Billion-scale similarity search with gpus. *arXiv* preprint arXiv:1702.08734.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526.
- Zhizhong Li and Derek Hoiem. 2017. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947.
- Chun Yang Liu, WenBo Sun, WenHan Chao, and Wanxiang Che. 2013. Convolution neural network for relation extraction. In *International Conference on Advanced Data Mining and Applications*, pages 231– 242. Springer.
- David Lopez-Paz and Marc'Aurelio Ranzato. 2017. Gradient episodic memory for continual learning. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 6467–6476.
  - Arun Mallya, Dillon Davis, and Svetlana Lazebnik. 2018. Piggyback: Adapting a single network to multiple tasks by learning to mask weights. In Proceedings of the European Conference on Computer Vision (ECCV), pages 67–82.
- Michael McCloskey and Neal J Cohen. 1989. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier.
- Mike Mintz, Steven Bills, Rion Snow, and Daniel Jurafsky. 2009. Distant supervision for relation extraction without labeled data. In Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP, pages 1003–1011, Suntec, Singapore. Association for Computational Linguistics.
- Makoto Miwa and Mohit Bansal. 2016. End-to-end relation extraction using LSTMs on sequences and tree

structures. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1105–1116, Berlin, Germany. Association for Computational Linguistics.

- Abiola Obamuyide and Andreas Vlachos. 2019. Metalearning improves lifelong relation extraction. In *Proceedings of the 4th Workshop on Representation Learning for NLP (RepL4NLP-2019)*, pages 224– 229, Florence, Italy. Association for Computational Linguistics.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Sachin Ravi and Hugo Larochelle. 2017. Optimization as a model for few-shot learning. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. 2017. icarl: Incremental classifier and representation learning. In 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017, pages 5533–5542. IEEE Computer Society.
- Haopeng Ren, Yi Cai, Xiaofeng Chen, Guohua Wang, and Qing Li. 2020. A two-phase prototypical network model for incremental few-shot relation classification. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 1618–1629, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Danilo Jimenez Rezende, Shakir Mohamed, Ivo Danihelka, Karol Gregor, and Daan Wierstra. 2016. One-shot generalization in deep generative models. In Proceedings of the 33nd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016, volume 48 of JMLR Workshop and Conference Proceedings, pages 1521–1529. JMLR.org.
- Hippolyt Ritter, Aleksandar Botev, and David Barber. 2018. Online structured laplace approximations for overcoming catastrophic forgetting. In Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pages 3742–3752.
- Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. 2016. Progressive neural networks. arXiv preprint arXiv:1606.04671.

992

993

994

995

996

997

998

999

946

947

948

949

950

- 892

- 899
- 900 901
- 902 903 904
- 905 906
- 907 908
- 909
- 910 911 912

913 914 915

- 916 917
- 919 921
- 923
- 924 925

927

928 929

- 930 931
- 932

935

937

938

939 940

941

944

- Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. 2016. Metalearning with memory-augmented neural networks. In International conference on machine learning, pages 1842-1850. PMLR.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1715-1725, Berlin, Germany. Association for Computational Linguistics.
- Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. 2017. Continual learning with deep generative replay. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 2990-2999.
  - Yusuke Shinyama and Satoshi Sekine. 2006. Preemptive information extraction using unrestricted relation discovery. In Proceedings of the Human Language Technology Conference of the NAACL, Main Conference, pages 304-311, New York City, USA. Association for Computational Linguistics.
- Livio Baldini Soares, Nicholas FitzGerald, Jeffrey Ling, and Tom Kwiatkowski. 2019. Matching the blanks: Distributional similarity for relation learning. arXiv preprint arXiv:1906.03158.
- Ang Sun, Ralph Grishman, and Satoshi Sekine. 2011. Semi-supervised relation extraction with large-scale word clustering. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, pages 521-529, Portland, Oregon, USA. Association for Computational Linguistics.
- Fan-Keng Sun, Cheng-Hao Ho, and Hung-Yi Lee. 2020. LAMOL: language modeling for lifelong language learning. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.
- Eleni Triantafillou, Richard S. Zemel, and Raquel Urtasun. 2017. Few-shot learning through an information retrieval lens. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 2255-2265.
- Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. Journal of machine learning research, 9(11).
- Hong Wang, Wenhan Xiong, Mo Yu, Xiaoxiao Guo, Shiyu Chang, and William Yang Wang. 2019. Sentence embedding alignment for lifelong relation extraction. 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

796–806, Minneapolis, Minnesota. Association for Computational Linguistics.

- Tongtong Wu, Xuekai Li, Yuan-Fang Li, Reza Haffari, Guilin Qi, Yujin Zhu, and Guoqiang Xu. 2021. Curriculum-meta learning for order-robust continual relation extraction. arXiv preprint arXiv:2101.01926.
- Hong-Ming Yang, Xu-Yao Zhang, Fei Yin, and Cheng-Lin Liu. 2018. Robust classification with convolutional prototype learning. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 3474–3482.
- Limin Yao, Aria Haghighi, Sebastian Riedel, and Andrew McCallum. 2011. Structured relation discovery using generative models. In Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, pages 1456–1466, Edinburgh, Scotland, UK. Association for Computational Linguistics.
- Mo Yu, Wenpeng Yin, Kazi Saidul Hasan, Cicero dos Santos, Bing Xiang, and Bowen Zhou. 2017. Improved neural relation detection for knowledge base question answering. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 571-581, Vancouver, Canada. Association for Computational Linguistics.
- Dmitry Zelenko, Chinatsu Aone, and Anthony Richardella. 2002. Kernel methods for relation extraction. In Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing (EMNLP 2002), pages 71-78. Association for Computational Linguistics.
- Daojian Zeng, Kang Liu, Yubo Chen, and Jun Zhao. 2015. Distant supervision for relation extraction via piecewise convolutional neural networks. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, pages 1753-1762, Lisbon, Portugal. Association for Computational Linguistics.
- Daojian Zeng, Kang Liu, Siwei Lai, Guangyou Zhou, and Jun Zhao. 2014. Relation classification via convolutional deep neural network. In Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers, pages 2335-2344, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.
- Friedemann Zenke, Ben Poole, and Surya Ganguli. 2017. Continual learning through synaptic intelligence. In Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017, volume 70 of Proceedings of Machine Learning Research, pages 3987-3995. PMLR.

- Chi Zhang, Nan Song, Guosheng Lin, Yun Zheng, Pan Pan, and Yinghui Xu. 2021. Few-shot incremental learning with continually evolved classifiers. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 12455– 12464.
  - Yuhao Zhang, Victor Zhong, Danqi Chen, Gabor Angeli, and Christopher D. Manning. 2017. Positionaware attention and supervised data improve slot filling. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 35–45, Copenhagen, Denmark. Association for Computational Linguistics.
    - Kai Zhu, Yang Cao, Wei Zhai, Jie Cheng, and Zheng-Jun Zha. 2021. Self-promoted prototype refinement for few-shot class-incremental learning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 6801–6810.

### A Appendix

1001

1002

1003

1004

1005

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1020

1021

### A.1 Difference between CRL and CFRL



Figure 5: Except for the first task which has enough training data, the subsequent new tasks are all *few-shot* in CFRL. In contrast, CRL assumes enough training data for every task.

### A.2 Block Diagram of ERDA Training



Figure 6: The block diagram of ERDA's training at time step k.

#### A.3 Hyperparameter Search

We follow EMAR (Han et al., 2020) and use a grid1023search to select the hyperparameters. Specifically,1024the search spaces are:1025

1022

1031

1032

1035

- Search range for  $\alpha$  is [0.3, 0.8] with a step size 1026 of 0.05.
- Search range for K is [1,3] with a step size of 1.
- Search range for  $m_1$  and  $m_2$  is [0.1, 0.3] with a step size of 0.1.
- Search range for  $m_3$  is [0.01, 0.03] with a step size of 0.01.
- Search range for  $iter_2$  in Alg. 1 is [1,3] with a step size of 1.





Figure 7: Relation extraction results for ERDA (our) and EMAR (Han et al., 2020) on the FewRel benchmark under the CRL setting. We randomly split the 80 relations into 8 tasks, where each task contains 10 relations. And we sample 100 examples per relation. From this figure, we can observe that ERDA outperforms EMAR in all CRL tasks.