Meta-learning with Prototypical Query Sample Editing for Low Resource Continual Learning

Anonymous ACL submission

Abstract

001 Memory replay, also known as experience 002 rehearsal, is a mainstream method to retain knowledge in continual learning (CL), especially in NLP. Recently, meta-learning has been introduced to augment memory replay. It serves to enable efficient knowledge transfer, thereby alleviating catastrophic forgetting. However, memory replay is often episodic to avoid over-fitting to past samples. Its sparse occurrence also limits the convergence of model on past samples, especially in low resource sce-011 012 narios. This paper aims to fully exploit the potential of meta-learning. We study the feasibility of solely using meta-learning to solve lifelong language learning. We propose an optimization-based meta-learning framework for CL in accordance with MAML. In particu-017 lar, we edit query information for meta-learning via a prototypical network. The meta-objective is modified to depict a CL scenario. We conduct extensive experiments on benchmark text 021 classification datasets. The results testify the superiority of our method in terms of forgetting mitigation, fast adaptation and memory efficiency in a low resource NLP scenario.

1 Introduction

026

027

028

041

Existing continual learning (CL) algorithms update model parameters whenever a new task arrives. The important information from earlier tasks can be easily erased or overwritten when the data distribution shifts. Consequently, catastrophic forgetting (McCloskey and Cohen, 1989) occurs and harms performance on preceding tasks. Especially for language models, the number of parameters in a language model is usually huge. A small change in parameter space could impact model output unexpectedly (Wang et al., 2019). Hence, replaying past examples is the mainstream approach for lifelong language learning. However, existing replay-based methods not only easily over-fit to past samples (Zhao et al., 2022), but also cause under-fitting problem on low resource tasks.

Meta-learning can enable efficient knowledge transfer. Recently, it has been introduced as a learning framework for CL. Mostly, meta-learning serves to perform efficient knowledge transfer via fast adaptation and facilitate memory replay. Memory replay is often episodic to avoid over-fitting problem. This episodic nature is unfavourable for cross-domain fast adaptation when training examples are insufficient and memory budget for past samples is small. 043

044

045

046

047

050

051

052

054

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

075

076

077

078

079

081

To address this dilemma, we propose a method to solve low resource CL problems, namely MAML-CL. Specifically, we use Model-Agnostic Meta-Learning (MAML) (Finn et al., 2017) as our metalearning framework. We argue that meta-learning can solve CL problems without episodic memory replay. The mainstream CL methods limit the change of non-trivial weights or parameters, hereby retaining knowledge. In meta-learning, the variation on parameter space is restricted by metaobjective. We use a prototypical network to select representative past examples for query set. It modifies the meta-objective and controls the learning process on the new task, thereby ameliorating catastrophic forgetting.

We conduct extensive experiments on CL benchmark datasets from (Zhang et al., 2015), popularized by (de Masson d'Autume et al., 2019) in lifelong language learning. This collection of datasets includes text classification datasets from diverse domains. We reduce the size of the training set for each task to its 10%, i.e., 11,500 per task. We demonstrate our method's robustness to forgetting mitigation, fast adaptation and memory efficiency in a low-resource scenario.

The contributions of this work are threefold: (1) Our method validates the feasibility of using meta-learning without sparse experience replay, the mainstream lifelong language learning approach; (2) We design a prototypical query example editing method to adapt the meta-objective to a CL scenario. Our method addresses low resource CL problems. It well-preserves fast adaptation ability of meta-learning while avoiding catastrophic forgetting; (3) We perform extensive experiments in a low resource scenario. The results testify the superiority of our method on forgetting mitigation, fast adaptation and memory efficiency.

2 Related Work

084

086

090

100

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

Robust continual learning aims to guarantee the *stability* of handling various tasks, while showing *plasticity* on a novel domain. Existing approaches can be categorised into two main mainstreams, i.e., regularization-based methods (Kirkpatrick et al., 2016; Li and Hoiem, 2016; Zenke et al., 2017) and replay-based methods (Wang et al., 2019, 2020; Chaudhry et al., 2019; Sun et al., 2019). In general, regularization-based methods retain the updates of non-trivial parameters or weights by adding constraints or penalty. Generally, language models are deep learning networks. Given the complexity of deep learning models, memory replay-based approaches are broadly considered as a plausible approach for lifelong language learning.

Replay-based methods save a certain amount of previous training data and retrieval them while training for a novel domain. Many models used random selection strategy to select or save past examples. But it always failed to consider sam-Recent replay-based methods ple efficiency. utilised episodic memory replay with a 1% replay rate. It was first proposed in MbPA++ (de Masson d'Autume et al., 2019). MbPA++ retrieved K nearest neighbours via Euclidean distance function for local adaptation, which is expensive in its inference process. To address this problem, (Holla et al., 2020) introduced a replay frequency to harness the replay sparseness in terms of time and size. Nevertheless, they neglected memory constraints. The size of occupied memory is explosively expanding with the increase of training samples.

Meta-learning in CL. Recently, meta-learning 124 has been introduced into CL models, considering 125 its ability of fast adaptation and knowledge transfer. 126 Recent works employed MAML (Finn et al., 2017) 127 to improve initial parameters. It can fast adapt to various domains with few learning samples. Meta-129 MbPA (Wang et al., 2020) performed local adapta-130 tion with episodic memory and used MAML to find 131 a better initialized state for local adaptation. OML-132 ER (Holla et al., 2020) and ANML-ER (Holla et al., 133

2020) utilised an online meta-learning model and a neuromodulated meta-learning respectively for fast adaptation, augmented with sparse experience replay. Additionally, some CL models used Reptile (Nichol et al., 2018) as their meta-learning algorithms. MER (Riemer et al., 2019) regularized the objective of experience replay via a modified Reptile (Nichol et al., 2018) algorithm and memory replay module. MLLRE (Obamuyide and Vlachos, 2019) also adopted Reptile to meta updates parameters via augmented training set.

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

160

161

162

163

164

165

166

167

168

169

170

171

172

3 Problem Formulation

3.1 Continual Learning

In a basic CL setup, the learner ingests a stream of training examples. We assume the training stream consists of K tasks, $\{T_1, T_2, ..., T_K\}$, in an ordered sequence. Given the latest task \mathcal{T}_K and a CL learner f_{θ_K} , the objectives are: (a) to finetune parameter set θ_K over a parameter space Θ , such that the task objective $\mathcal{L}_{\mathcal{T}_K}$ is minimal. That is,

$$\tilde{\theta}_K = \arg\min_{\theta_K \in \Theta} \mathcal{L}_{\mathcal{T}_K}(\theta_K) \tag{1}$$

(b) to perform well with the learned parameter set $\tilde{\theta}_K$ on all preceding tasks $\{\mathcal{T}_1, \mathcal{T}_2, ..., \mathcal{T}_{K-1}\}$ without using the past training sets. Assuming that all tasks are equally important, the overall objective is to minimise the expected risk of all seen tasks, $\{\mathcal{T}_1, \mathcal{T}_2, ..., \mathcal{T}_K\}$, as:

$$\min_{\tilde{\theta}_K} \sum_{k=1}^K \mathbb{E}_{\mathcal{T}_k}[\mathcal{L}_{\mathcal{T}_k}(\tilde{\theta}_K)]$$
(2)

A replay-based method allows the model to retain a certain amount of past instances. Typically, the basic CL setup allows the memory buffer size to a constant B, where the amount of past training data in memory should not exceed B.

In this paper, we consider a popular scenario of continual learning, i.e., *class-incremental learning* (CIL), where task identity information is not provided in inference, and a single classifier is for all classes. We leverage cross entropy loss \mathcal{L}_{CE} as the loss of a task \mathcal{T}_k as,

$$\mathcal{L}_{\mathcal{T}_{k}}(\theta) = \mathcal{L}_{CE}(\boldsymbol{x}, \boldsymbol{y}; \theta)$$
$$= -\sum_{j=1}^{|D_{k}|} \sum_{i=1}^{N_{k}} y_{ji} \log(\sigma(f_{\theta}(x_{j}))_{i})$$
(3) (3)

where k is the task index, $D_k \subseteq D_{\text{train}}^k, D_{\text{train}}^k$ is 174 the ground truth label set of \mathcal{T}_k, N_k is the number 175

216 217 218

219

.

221

222

225

of classes in \mathcal{T}_k and σ is the activation function (e.g., sigmoid or softmax).

3.2 Meta-learning

176

177

178

179

180

181

184 185

187

188

189

190

191

193

194

195

196

197

198

199

201

207

209

211

212

213

Model agnostic meta-learning (MAML) (Finn et al., 2017) algorithm is known as "think in advance" by learning an optimal initial state of an algorithm. We assume a task T and a set of initial parameters ϕ over a parameters space Φ . We expect ϕ to yield a minimal loss after m updates in Tas,

$$\min_{\phi} \mathbb{E}_{\mathcal{T}}[\mathcal{L}_{\mathcal{T}}(U_{\mathcal{T}}^m(\phi))] \tag{4}$$

where $U_{\mathcal{T}}^m$ is the update operation that performs m times gradient-based updates on parameters ϕ , using samples drawn from distribution $p(\mathcal{T})$. In an online MAML setting, each episode uses m batches of training instances as a support set S. The inner loop algorithm performs m-step SGD on parameters ϕ in the inner loop, as:

$$\phi^* = U_{\mathcal{T}}^m(\phi)$$

= $U_{\mathcal{T},S}(\phi)$ (5)
= $\phi - \alpha \nabla_{\phi} \mathcal{L}_{\mathcal{T}}(S,\phi)$

where α is the stepsize. Test examples that specified problems in outer loop algorithm are a set of query samples Q. The updated ϕ^* are further optimized on Q to achieve meta objective as,

$$\min_{\phi} \mathbb{E}_{\mathcal{T}}[\mathcal{L}_{\mathcal{T},Q}(\phi^*)] = \min_{\phi} \mathbb{E}_{\mathcal{T}}[\mathcal{L}_{\mathcal{T},Q}(U_{\mathcal{T},S}(\phi))]$$
(6)

4 Prototypical Query Sample Editing

To align meta-objective with CL objective, we edit query set Q to generalize all seen tasks, $\{\mathcal{T}_1, \mathcal{T}_2, ..., \mathcal{T}_K\}$. Particularly, we select representative past examples as query examples for metalearning. We apply a prototypical network for query sample selection, where each class is characterized by a prototype.

4.1 Dynamic Prototype Computation

We devise a prototypical network, $h_{\phi_{\text{proto}}}$, for prototype computation. Each prototype c_l is computed as the mean vector of the embedded support examples from the same class l, as:

$$\mathbf{c}_l \leftarrow \frac{1}{|S_l|} \sum_{(x_j, y_j) \in S_l} h_{\phi_{proto}}(x_j) \tag{7}$$

where l is the class index and S_l is the support set for class l. The devised prototypical network $h_{\phi_{\text{proto}}}$ is meta-learnable and dynamically compute and update prototypes in each episode. And, we include the prototypical loss J_P in inner loop loss $\mathcal{L}_{\text{inner}}$ as a regularization term. For a task \mathcal{T}_k , the prototypical loss J_P (Snell et al., 2017) is,

$$J_P(\phi_{\text{proto}}) = \frac{1}{N_k} \sum_{l=1}^{N_k} (\frac{1}{|Q_l|} [d(h_{\phi_{\text{proto}}}(\boldsymbol{x}_l), \mathbf{c}_l) + \log \sum_{l'} \exp(-d(h_{\phi_{\text{proto}}}(\boldsymbol{x}_l), \mathbf{c}'_l))])$$
(8)

where Q_l is the query set for class l, $(\boldsymbol{x}_l, y_l) \in Q_l$, N_k is the number of classes in \mathcal{T}_k and $d(\cdot)$ is a distance function. In this paper, we leverage Euclidean distance.

A	lgorithm 1: Meta-training
	Input: Initial parameters $\theta = \phi_{\text{proto}} \cup \phi_{\text{pred}}$,
	training set D_{train} , support set size m ,
	episode index i , memory buffer D_{Memory} , No.
	of support examples for each class N_S , No. of
	query examples for each class N_Q , inner-loop
	learning rate α , outer-loop learning rate β ,
	and No. of selected samples per class N_{select} .
	Output: Trained parameters θ and memory D_{Memory}
1	for $i = 1, 2,$ do
2	$S \leftarrow m$ batches of examples D_{train}^i from the
	stream
3	[Prototype Computation]
4	$J_p(\phi_{\text{proto}}) \leftarrow 0$
5	for class l in S do
6	Select N_S examples of class l from S as S_l .
7	Select N_Q examples of class l from $S \setminus S_l$ as Q_l .
8	Compute and update prototype c_l as Eq.7
	and update c_l in prototype set c_l .
9	Compute and update prototypical loss
	$J_n(\phi_{\text{proto}})$ as Eqn.8.
10	end
11	[Inner Loop]
12	Perform SGD on ϕ_{pred} to minimize Eqn.10 as
	$\phi_{\text{pred}}' = \phi_{\text{pred}} - \alpha \nabla_{\theta} \mathcal{L}_{inner}(S; \theta).$
13	[Prototypical Query Sample Editing]
14	for class l in prototype set c do
15	Select N_{select} nearest examples to c_l from
	$S \cup D_{ ext{Memory}}.$
16	Update examples of class l in D_{Memory} .
17	end
18	$Q \leftarrow D_{\text{Memory}}$ OR randomly sample a batch of
	data from D_{Memory}
19	[Outer Loop]
20	Perform Adam update on θ' to minimize Eqn.11
	as $\theta \leftarrow \theta - \beta \nabla_{\theta'} \mathcal{L}_{outer}(Q; \theta')$, where
	$ heta' = \phi_{ m proto} \cup \phi'_{ m pred}.$
21	if all training data are seen then
22	Stop Iteration
23	end
24	end

4.2 Query Sample Selection

The synthetic prototypes serve as a reference for query sample selection. In each episode, we choose query samples from training batches and memory set D_{memory} . In particular, we choose top N_{select} instances for each class with the highest similarity score to their corresponding prototype. The similarity score is calculated via Euclidean distance. The newly selected samples replace past ones and update memory set. The query samples of \mathcal{T}_k is denoted as Q_k , in which we expect that $(\boldsymbol{x}, \boldsymbol{y}) \in Q_k \sim p(\mathcal{T}_k)$. We propose two read functions to obtain query samples: (a) read *all* from memory; (b) read *randomly* from memory, where the amount of retrieval instances equals to the minibatch size, *b*.

5 Meta-learning with Prototypical Query Sample Editing

The meta-objective in outer loop governs the optimization process of task-specific learning in inner loop. We edit query samples by a prototypical network (Snell et al., 2017) to revise meta-objective to CL objective. We control the variation on parameter space by the edited meta-objective, hereby solving CL problems.

5.1 Edited Meta-objective for Continual Learning

We edit query examples via prototypes for outer loop algorithm. Then, we formulate a new metaobjective edited by prototypical information as:

$$\min_{\theta_{K}} \mathbb{E}_{\mathcal{T}}[\mathcal{L}_{\mathcal{T},Q}(U_{\mathcal{T}_{K},S}(\theta_{K}))]$$

$$= \min_{\theta_{K}} \mathbb{E}_{\mathcal{T}}[\mathcal{L}_{\mathcal{T},Q}(\tilde{\theta}_{K})]$$

$$= \min_{\theta_{K}} \sum_{k=1}^{K} \mathbb{E}_{\mathcal{T}_{k}}[\mathcal{L}_{\mathcal{T}_{k},Q_{k}}(\tilde{\theta}_{K})]$$

$$= \min_{\theta_{K}} \sum_{k=1}^{K} \mathbb{E}_{(\boldsymbol{x},\boldsymbol{y})\sim p(\mathcal{T}_{k})}[\mathcal{L}_{CE}(\boldsymbol{x},\boldsymbol{y};\tilde{\theta}_{K})]$$

$$= \min_{\theta_{K}} \sum_{k=1}^{K} \mathbb{E}_{\mathcal{T}_{k}}[\mathcal{L}_{\mathcal{T}_{k}}(\tilde{\theta}_{K})]$$
(9)

Evidently, the edited meta-objective has the same
expected loss as that of CL as shown in Eqn.2, but
different optimization parameters.

260 Prevent catastrophic forgetting: The consis261 tent expected loss indicates that the meta-learning
262 framework can depict a CL scenario. Given a

learned parameter set $\tilde{\theta}_K$ heavily biased towards the distribution of the latest task, $p(\mathcal{T}_K)$, the objective is to perform well on all tasks, $\{\mathcal{T}_1, \mathcal{T}_2, ..., \mathcal{T}_K\}$. Hence, we interpret CL problems into a form of meta-learning.

263

264

265

267

268

269

270

271

272

273

274

275

276

277

280

281

283

284

285

287

290

291

292

293

Preserve fast adaptation: Different optimization parameters suggest different learning strategy. As shown in Eqn.2, the optimization parameters of CL objective are finetuned parameters $\tilde{\theta}_K$. It neglects the next optimal steps, potentially causing more update steps for learning different domains. Conversely, the optimization parameters of edited meta-objective are initial parameters θ_K . It learns an optimal initialization state and yields a minimal loss with few update steps on all seen tasks. Hence, it preserves fast adaptation ability of meta-learning while solving catastrophic forgetting.

6 Model

The proposed model f_{θ} consists of a representation learning network $h_{\phi_{\text{proto}}}$ with learnable parameters ϕ_{proto} and a prediction network $g_{\phi_{\text{pred}}}$ with learnable parameters ϕ_{pred} . It is described as $f_{\theta}(x) = g_{\phi_{\text{pred}}}(h_{\phi_{\text{proto}}}(x))$. The prototypical network is a single-hidden-layer feed-forward neural network on top of an encoder. The prediction network is a single linear layer followed by a softmax.

Meta-training. We sample *m* batches of examples instantaneously from the data stream D_{train} as support set *S* for inner loop. The inner loop loss $\mathcal{L}_{\text{inner}}$ includes the current task loss $\mathcal{L}_{\mathcal{T}_K}$, i.e., cross entropy loss \mathcal{L}_{CE} , and prototypical network loss J_P ,

$$\mathcal{L}_{\text{inner}}(S;\theta) = \mathcal{L}_{CE}(S;\theta) + J_P(\phi_{\text{proto}})$$

= $\mathcal{L}_{CE}(S;\phi_{\text{proto}} \cup \phi_{\text{pred}}) + J_P(\phi_{\text{proto}})$ 295
(10)

The inner loop algorithm only finetunes PN model 296 $g_{\phi_{\text{pred}}}$ via SGD with learning rate α , $\phi'_{\text{pred}} =$ 297 $\phi_{\text{pred}} - \alpha \nabla_{\theta} \mathcal{L}_{\text{inner}}(S; \theta)$. In outer loop algorithm, 298 both parameters ϕ_{proto} and ϕ'_{pred} are meta-trained 299 on query set Q. The outer loop loss is, 300

$$\mathcal{L}_{\text{outer}}(Q; \theta') = \sum_{k=1}^{K} \mathbb{E}_{\mathcal{T}_{k}}[\mathcal{L}_{\mathcal{T}_{k}}(\theta')]$$
$$= \sum_{k=1}^{K} \mathbb{E}_{(\boldsymbol{x}, \boldsymbol{y}) \sim p(\mathcal{T}_{k})}[\mathcal{L}_{CE}(\boldsymbol{x}, \boldsymbol{y}; \theta')]$$
$$= \mathcal{L}_{CE}(Q; \theta')$$
(11)

256

226

228

229

230

234

235

237

240

241

242

243

247

251

252

253

where $heta' = \phi_{
m proto} \cup \phi'_{
m pred}$. The optimizer of 302 the outer loop algorithm is Adam optimizer with 303 a learning rate β . The optimization process is 304 $\theta \leftarrow \theta - \beta \nabla_{\theta} \mathcal{L}_{outer}(Q; \theta')$, where the gradients are computed with respect to the initial parameters θ . To reduce computationally expensive, 307 we use the first-order approximation, namely FO-MAML (Finn et al., 2017). The gradients are computed with respect to the finetuned parame-310 ters θ' . Then, the outer loop optimization process 311 is $\theta \leftarrow \theta - \beta \nabla_{\theta'} \mathcal{L}_{outer}(Q; \theta')$. The meta-training 312 process is shown in Algorithm 1. 313

Meta-inference. We randomly sample *m* batches of examples from memory set D_{Memory} as the support set *S* for inner loop. We perform SGD on ϕ_{pred} to minimize cross entropy loss $\mathcal{L}_{CE}(S;\theta)$ as $\phi'_{\text{pred}} = \phi_{\text{pred}} - \alpha \nabla_{\theta} \mathcal{L}_{\text{inner}}(S;\theta)$. Then, we output the predication on the test samples $\boldsymbol{x}_{\text{test}}$ as $\hat{\boldsymbol{y}}_{\text{test}} = f_{\theta}(\boldsymbol{x}_{\text{test}}) = g_{\phi'_{\text{pred}}}(h_{\phi_{\text{proto}}}(\boldsymbol{x}_{\text{test}}))$.

7 Experiments

7.1 Datasets

314

315

316

317

318

320

321

326

327

331

332

333

334

We use the collection of text classification datasets from (Zhang et al., 2015)¹, including AGNews (news classification; 4 classes), Yelp (sentiment analysis; 5 classes), Amazon (sentiment analysis; 5 classes), DBpedia (Wikipedia article classification; 14 classes) and Yahoo (questions and answers categorization; 10 classes). Following prior work (de Masson d'Autume et al., 2019; Holla et al., 2020), we use the balanced version of the collection and merge the classes of Yelp and Amazon. Thus, we have 33 classes in total. We randomly sample 115,000 training examples and 7,600 test examples from each of the datasets. Each dataset is seen as a separate learning task.

Table 1: Input Datasets Orders

No.	Dataset Orders
1	$Yelp \rightarrow AGNews \rightarrow Amazon$
2	$Yelp \rightarrow Amazon \rightarrow AGNews$
3	Amazon \rightarrow Yelp \rightarrow AGNews
4	Amazon \rightarrow AGNews \rightarrow Yelp
5	$AGNews \rightarrow Yelp \rightarrow Amazon$
6	$AGNews \rightarrow Amazon \rightarrow Yelp$
7	$Yelp \rightarrow AGNews \rightarrow DBpedia \rightarrow Amazon \rightarrow Yahoo$
8	$DBpedia \rightarrow Yahoo \rightarrow AGNews \rightarrow Amazon \rightarrow Yelp$
9	$Yelp \rightarrow Yahoo \rightarrow Amazon \rightarrow DBpedia \rightarrow AGNews$
10	$AGNews \rightarrow Yelp \rightarrow Amazon \rightarrow Yahoo \rightarrow DBpedia$

¹http://goo.gl/JyCnZq

7.2 Baselines

• MAML-SEQ is a sequential learning model using the same meta-learning (Finn et al., 2017) framework as ours without memory replay. 337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

354

355

356

357

359

360

361

362

363

364

365

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

383

- **Replay** is a sequential learning model with episodic experience replay. It applies the common-used write and read mechanisms, i.e., write all and randomly read.
- **A-GEM** (Chaudhry et al., 2019) with episodic memory replay is a commonly-used lifelong language learning baseline, where catastrophic forgetting is prevented by gradient constraints.
- **OML-ER** (Holla et al., 2020) uses the same online meta-learning framework as ours but with episodic memory replay. Its default setting is writing all and randomly read.
- **PMR** (Ho et al., 2021) uses the same prototypical sample selection criteria as ours but use episodic memory replay.

7.3 Experimental Setup

We define our low resource scenario as insufficient training examples and strict memory constraints. We reduce the size of the training set to its 10%, i.e., 11,500 per task. All models are trained in one epoch. We consider the effect of different orders of datasets. We conduct experiments on all six possible orders of three datasets, i.e., Yelp, AGNews and Amazon, as shown in Table 1. Following prior work (de Masson d'Autume et al., 2019; Holla et al., 2020), we also run experiments on all 5 tasks, using 4 possible orders.

The example encoder is a pre-trained ALBERT-Base-v2 (12M) from Hugging Face Transformers (Lan et al., 2020), where we truncate the input sequence length to 200. For meta-learning models, we use SGD as their inner loop optimizer with a learning rate, $\alpha = 1e^{-3}$, and Adam as their outer loop optimizer with a learning rate, $\beta = 3e^{-5}$. For A-GEM (Chaudhry et al., 2019) and Replay, we use Adam as their optimizer with learning rate, $3e^{-5}$. The number of mini-batches in each epoch, m =5. Models with prototypical network samples data uniformly at random without replacement, where batch sizes depend on the number of classes from the training sets, e.g., b = 20 for AGNews and b

Mathad	Memory		Average					
Method	Constraint (B)	(1)	(2)	(3)	(4)	(5)	(6)	Accuarcy
MAML-SEQ	No data saved	34.6	38.8	32.4	29.3	22.6	28.0	30.9 ± 5.6
AGEM		38.4	29.8	29.9	40.4	37.3	37.8	35.6 ± 4.6
Replay	Unlimited 100% seen data	42.5	29.8	29.6	46.5	40.4	38.7	37.9 ± 6.9
OML-ER	Chilinited, 100 % seen data	47.8	25.4	31.4	40.8	48.5	39.2	38.9 ± 9.1
PMR		26.9	26.4	21.6	26.3	19.8	25.2	24.4 ± 3.0
Ours.(all)	5 per class, 0.13 % seen data	56.3	53.6	49.4	54.2	47.9	45.6	51.2 ± 4.2
Ours.(random)		52.1	47.1	41.4	49.0	43.0	39.1	45.3 ± 4.9

Table 2: Accuracy on 3 Tasks with 6 Different Orderings

Order Index	Mathod	Yelp		AGNews		Amazon		Overall Forgetting		
Order maex	Wiethou	$A_{\rm CL,1}$	F_1	$A_{\rm CL,2}$	F_2	$A_{\rm CL,3}$	F_3	Overall Forgetting		
	OML-ER	46.8	3.26	49.2	36.8	47.3	3.2	43.3		
Order (1)	Ours.(all)	44.7	-1.3	79.5	2.8	44.5	-4.9	-3.4		
	Ours.(random)	41.3	0.4	74.3	6.4	40.7	-4.5	2.3		
Order Index	Method	AGNews		Yelp		Amazon		Overall Forgetting		
Order maex	Method	$A_{\rm CL,1}$	F_1	$A_{\rm CL,2}$	F_2	$A_{\rm CL,3}$	F_3	Overall Forgetting		
	OML-ER	58.4	27.7	43.4	6.8	43.9	6.6	41.2		
Order (5)	Ours.(all)	68.0	14.3	38.2	5.2	37.6	2.0	21.5		
	Ours.(random)	74.1	6.6	26.9	14.8	28.0	8.2	29.6		

Table 3: Per-task and Overall Forgetting on 3 Tasks

= 25 for Yelp and Amazon. We set $N_S = 3 * m$ and $N_Q = 2 * m$ for prototypical loss computation. Other models utilise the random sampler, which randomly samples examples with batch size, b =25. The write rate p_{write} is set to 1. All models are executed on a Linux platform with 8 Nvidia Tesla A100 GPU and 40 GB of RAM. All experiments are performed using PyTorch (Paszke et al., 2019).

7.4 Evaluation Metrics

We evaluate the models after learning all tasks. The evaluation metrics are accuracy and forgetting. Let $A_{CL,k}$ be the macro-averaged accuracy on \mathcal{T}_k , the overall accuracy after learning a sequence of K tasks is ACC = $\frac{1}{K} \sum_{k=1}^{K} A_{CL,k}$. Let $A_{\text{single},k}$ be the accuracy of learning one task \mathcal{T}_k . The forgetting on a sequence of K tasks is $F = \sum_{k=1}^{K} F_k = \sum_{k=1}^{K} A_{\text{single},k} - A_{\text{CL},k}$, where F_k is the forgetting on task \mathcal{T}_k and describes how the updated parameters affect the performance on \mathcal{T}_k after sequential learning all tasks.

7.5 Results

394

400

401

402

403

404

405

406

407

408

409

410

411

We run experiment on all 6 possible orders of training sets while learning a sequence of 3 tasks. We record the average of 3 best results from 5 trials.

Accuracy. As shown in Table 2, our model with read all function, obtains the best results in most of the training sequences, except for Order (5), which is only 0.6% smaller than the best result. Given downsized training sets, it maintains average accuracy to more than 50%, suggesting its fast adaptation ability. Its standard deviations of the average accuracy is 4.2%, which is smaller than most baselines. It shows a good robustness in terms of training set orders. The proposed method with read random function has the second-highest average accuracy. Under the same memory constraint, PMR is inferior to our models by at least 20%. Remarkably, both of the proposed models even outperform models with unlimited memory budget. The results demonstrate the superiority of our model in a low resource scenario.

412

413

414

415

416

417

418

419

420

421

422

423

424

426

427

428

429

430

431

432

434

435

436

437

439

440

441

Forgetting Measurements. Evaluation on for-425 getting shows model's ability of knowledge retention, but not directly indicating model performance. A negative value of forgetting suggests accuracy improvements. As shown in Table 2, our models obtain the best performance on Order (1) and the worst performance on Order (5). We compare our models with all baselines on these two orders (see Table 8 in Appendix). We show the result of our 433 models and a strong baseline, OML-ER, in Table 3. For Order (1), both of our models have the negative score of forgetting on Amazon. It demonstrates a positive forward knowledge transfer occurs. All models obtain a relatively high forgetting score on 438 AGNews. It is reasonable considering its distribution is quite different from Yelp and Amazon. As for Order (5), the last two tasks, Yelp and Ama-

Memory Constraint (B)	Method	Yelp	AGNews	Amazon	Overall Accuracy
B = 27	OML-ER	49.8	19.4	49.1	39.4
D = 97 (2 par along)	Ours.(all)	44.0	67.6	43.8	51.8
D = 27 (3 per class)	Ours.(random)	36.2	63.6	35.1	44.9
B = 45	OML-ER	48.0	13.1	49.3	36.8
P = 4F (5 man alogg)	Ours.(all)	44.7	79.5	44.5	56.3
D = 45 (5 per class)	Ours.(random)	41.3	74.3	40.7	52.1
B = 63	OML-ER	52.7	24.6	50.3	42.5
B = 63 (7 per class)	Ours.(all)	55.1	86.2	52.2	64.6
D = 05 (7 per class)	Ours.(random)	40.9	69.3	39.3	49.8
	OML-ER	46.8	49.2	47.3	47.8
B = 34,500 (Unlimited)	Ours.(all)	_	_	_	_
	Ours.(random)	55.1	86.2	52.2	64.6

Table 4: Per-task and Overall Accuracy with Different Memory Constraints

 Table 5: Comparison of Various Query Sample Selection Methods

Query Sample Selection	Yelp	AGNews	Amazon	Overall Accuracy	
Random		43.6	69.4	43.7	52.2
Prototypical Sample Selection	Diversity (Ours.)	44.7	79.5	44.5	56.3
Prototypical Sample Selection	Uncertainty	28.2	47.8	27.8	34.6

zon, are both from product reviews. As a result, 442 OML-ER has enough iterations to perform mem-443 ory replay on AGNews. Its accuracy on AGNews 444 445 improves by 9.1%. our methods use few samples to represent prior task, the increase of iterations make 446 them converge too well to generalize. Thereby, our 447 model with read all function suffers over-fitting 448 problems, especially on AGNews. Read random 449 function lessens this problem on AGNews, but de-450 grades the performance on Yelp and Amazon. We 451 consider it as the trade-off between representation 452 453 and generalization. Still, both of our models outperform OML-ER and address catastrophic forgetting. 454

7.6 Further Analysis

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

We use Order (1) to perform further analysis. Note that Yelp and Amazon datasets are product reviews, while AGNews are news. The data distributions are different. Hence, we pay extra attentions to model performance on AGNews.

Memory Efficiency. We compare our models with all baselines (see Table 9 in the Appendix). We show the result of our models and a strong baseline, OML-ER, in Table 4. Our method with read *all*² function achieves highest overall performance given different memory constraints. Surprisingly, its results are even better than OML-ER without memory constraint by at most 16.8%. Our method with read *random* function is not as good as our read *all* method. It suggests that read *all* function is more beneficial when giving a strict memory limitation. The result indicates memory efficiency of the proposed method. Also, both of our models outperform OML-ER on AGNews by a large margin. It shows our models' ability of knowledge retention.

470

471

472

473

474

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

Effect of Query Sample Selection. Table 5 compares different query sample selection strategies. We consider two main strategies, i.e., random and prototypical sample selections. We further study two popular paradigms in active learning, i.e., diversity-based and uncertainty-based method. Our proposed method is considered as a diversity-based method. While, opting for examples far from prototypes is an uncertainty-based method. The read function is read *all*. In general, random selection is seen as an efficient selection criteria. Diversity-based method outperforms random selection by more than 4% and outperforms uncertainty-based method by more than 20%.

Effect of Query Sample Editing. We compare our approach to the mainstream CL method, i.e., episodic memory replay. The sample selection method is prototypical sample selection and write function writes all selected samples. Table 6 shows our approach surpasses episodic memory replay in two different read functions. For read *all* function, it is superior to episodic memory replay by nearly 20%. The performance on AGNews exceeds episodic memory replay by more than 44%. The reason behind this bad performance of episodic

²Ours.(*all*) with unlimited memory means writing all data to memory and replaying all data, which is not possible for large scale dataset and can result in out-of-memory errors.

Continual Leaning Strategy	Read Function	Yelp	AGNews	Amazon	Overall Accuracy
Episodic Memory Replay	A 11	34.8	34.9	34.2	34.6
Ours.(all)	All	44.7	79.5	44.5	56.3
Episodic Memory Replay	Dandom	32.8	19.5	31.6	27.8
Ours.(random)	Kalluolli	41.3	74.3	40.7	52.1

Table 6: Comparison Between Episodic Memory Replay and Our Method

Mathad	No. of Training	Mamory Constraint		Order	Index		Average Accuracy	
Method	Samples	Memory Constraint	(5)	(6)	(7)	(8)	Average Accuracy	
A-GEM [†]			70.7	65.9	67.5	63.6	66.9 ± 3.0	
$MbPA++^{\dagger}$			70.8	70.9	70.2	70.7	70.7 ± 0.3	
REPLAY [‡]		575,000 (Unlimited)	69.5	66.2	65.2	68.3	67.3 ± 2.0	
OML-ER [‡]	575.000		75.4	76.5	75.4	75.4	75.7 ± 0.6	
PMR	575,000		61.2	65.7	66.1	55.9	62.2 ± 4.8	
Ours.(all)		165 (0.03 % seen data)	54.2	64.1	64.6	48.1	57.8 ± 8.0	
Ours.(random)			59.2	61.6	64.6	51.3	59.2 ± 5.7	
A-GEM			24.7	33.2	25.6	22.8	26.6 ± 4.7	
Replay		57,500 (Unlimited)	31.0	50.8	38.8	37.2	39.5 ± 8.3	
OML-ER	57 500		43.7	53.3	44.2	39.4	45.2 ± 5.8	
PMR	57,500		2.6	45.2	5.0	5.7	14.6 ± 20.4	
Ours.(all)		165 (0.03 % seen data)	56.8	60.7	65.2	51.3	58.5 ± 5.9	
Ours.(random)			50.6	59.5	48.7	36.1	48.7 ± 9.6	

Table 7: Accuracy on 5 Tasks with 4 Different Orderings

memory replay is its insufficient training iterations.
Consequently, it is prone to forgetting. Similarly,
for read *random* function, our method outperforms
episodic memory replay by a large margin. In addition, high replay frequency often leads to overfitting problem, especially when the amount of
memory samples is small. Our approach revisits a
small amount of past samples every iteration. But
it still maintains a more than 52% overall accuracy.

Learning More Tasks and Examples. We run 510 experiments on 5 datasets in 4 different orderings. 511 For comparison, the experimental setup follows 512 prior work (de Masson d'Autume et al., 2019)³, 513 514 where encoder is a pretrained BERT (Devlin et al., 2019). As shown in Table 7, given full training 515 examples and strict memory constraint, both of our 516 models are inferior to PMR by $3.0 \sim 4.4\%$. The over-fitting problem occurs when the training it-518 erations increases. Hence, our method with read 519 random function performs better than that with read 520 all function. When we downsize the training sets 522 to their 10%, our models are superior, validating their fast adaptation abilities. Furthermore, PMR 523 performs surprisingly well on Order (6), compared 524 to its performance on other training set orders. The 525 last two tasks in Order (6) are Yelp and Amazon. 526 527 These two datasets are both product reviews and

from the same domain. It increases the number of training iterations for PMR to learn different probability distributions via episodic memory replay. Thus, it results in a better accuracy.

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

8 Conclusion

In this paper, we propose an enhanced metalearning framework for continual text classification model learning, where we use a prototypical network to edit query samples and adapt metaobjective for continual learning. The experimental results manifest that our method has an outstanding performance in a low resource scenario, i.e., insufficient training examples and strict memory constraints. It also validates our method ensures fast adaptation while preventing catastrophic forgetting.

9 Limitations

Our method works well in a low resource scenario, where the amount of training data is not plentiful and the memory budget is strictly limited. Its performance also relies heavily on pretrained language models. In addition, the metalearning framework we used, namely MAML, is the standard framework. The effect of different meta-learning frameworks should be studied. We leave this investigation to future work. Furthermore, we can extend our model to other NLP tasks.

 $^{^{3\}dagger}$ Results obtained from (de Masson d'Autume et al., 2019) and [‡] Results obtained from (Holla et al., 2020).

References

- 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.
- Cyprien de Masson d'Autume, Sebastian Ruder, Lingpeng Kong, and Dani Yogatama. 2019. Episodic memory in lifelong language learning. In Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pages 13122–13131.
- 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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 4171–4186. Association for Computational Linguistics.
- 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.
- Stella Ho, Ming Liu, Lan Du, Longxiang Gao, and Yong Xiang. 2021. Prototypes-guided memory replay for continual learning. *CoRR*, abs/2108.12641.
- Nithin Holla, Pushkar Mishra, Helen Yannakoudakis, and Ekaterina Shutova. 2020. Meta-learning with sparse experience replay for lifelong language learning. *CoRR*, abs/2009.04891.
- James Kirkpatrick, Razvan Pascanu, Neil C. Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. 2016. Overcoming catastrophic forgetting in neural networks. *CoRR*, abs/1612.00796.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut.
 2020. ALBERT: A lite BERT for self-supervised learning of language representations. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.
- Zhizhong Li and Derek Hoiem. 2016. Learning without forgetting. In Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV, volume

9908 of *Lecture Notes in Computer Science*, pages 614–629. Springer.

- M. McCloskey and N. J. Cohen. 1989. Catastrophic interference in connectionist networks: The sequential learning problem. *Psychology of Learning and Motivation*, 24:109–165.
- Alex Nichol, Joshua Achiam, and John Schulman. 2018. On first-order meta-learning algorithms. *CoRR*, abs/1803.02999.
- 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.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems* 32, pages 8024–8035. Curran Associates, Inc.
- Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauro. 2019. Learning to learn without forgetting by maximizing transfer and minimizing interference. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net.
- Jake Snell, Kevin Swersky, and Richard S. Zemel. 2017. Prototypical networks for few-shot 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 4077–4087.
- Fan-Keng Sun, Cheng-Hao Ho, and Hung-Yi Lee. 2019. LAMAL: language modeling is all you need for lifelong language learning. *CoRR*, abs/1909.03329.
- 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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 796–806. Association for Computational Linguistics.
- Zirui Wang, Sanket Vaibhav Mehta, Barnabás Póczos, and Jaime G. Carbonell. 2020. Efficient meta lifelong-learning with limited memory. In *Proceedings of the 2020 Conference on Empirical Methods in*

667

555

560

561

565

566

567

568

569

571

572

573

574

576

577

579

585

589

590

593

594

595

596

598

599

609

610

- *Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 535–548. 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.
 - Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada, pages 649–657.
 - Kang Zhao, Hua Xu, Jiangong Yang, and Kai Gao. 2022. Consistent representation learning for continual relation extraction. In *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 3402–3411. Association for Computational Linguistics.

A Appendix

671 672

673

674

675

676

677

678

679 680

681

682

683 684

685

686

687

688

689

690

Order Index	Method	Yelp		AGNews		Amazon		Overall Forgetting
Order maex	Method	$A_{\rm CL,1}$	F_1	$A_{\rm CL,2}$	F_2	$A_{\rm CL,3}$	F_3	Overall Polgetting
	MAML-SEQ	38.5	10.7	32.7	51.3	32.6	8.4	70.4
	AGEM	54.9	3.6	4.7	85.3	55.4	1.8	90.7
	Replay	56.7	1.8	13.3	76.7	57.5	-0.3	78.2
Order (1)	OML-ER	46.8	3.26	49.2	36.8	47.3	3.2	43.3
	PMR	41.6	1.6	0.2	84.9	38.9	-0.9	85.6
	Ours.(all)	44.7	-1.3	79.5	2.8	44.5	-4.9	-3.4
	Ours.(random)	41.3	0.4	74.3	6.4	40.7	-4.5	2.3
	Mathad							
Order Index	Method	AGN	ews	Yel	р	Ama	zon	Overall Forgetting
Order Index	Method	AGN $A_{CL,1}$	ews F1	Yel $A_{\rm CL,2}$	р <i>F</i> 2	Ama $A_{CL,3}$	zon F ₃	Overall Forgetting
Order Index	Method MAML-SEQ	AGN A _{CL,1} 27.9	ews F1 56.1	Yel A _{CL,2} 20.4	p <u>F2</u> 28.8	Ama: A _{CL,3} 19.4	<i>F</i> 3 21.6	Overall Forgetting 106.5
Order Index	Method MAML-SEQ AGEM	AGN A _{CL,1} 27.9 0.1	$ ews F_1 56.1 89.9 $	Yel A _{CL,2} 20.4 55.8		Ama A _{CL,3} 19.4 56.1	F_3 21.6 1.1	Overall Forgetting 106.5 93.7
Order Index	Method MAML-SEQ AGEM Replay	AGN A _{CL,1} 27.9 0.1 7.6	ews F_1 56.1 89.9 82.4	Yel A _{CL,2} 20.4 55.8 56.5		Ama A _{CL,3} 19.4 56.1 57.2	zon F_3 21.6 1.1 0.0	Overall Forgetting 106.5 93.7 84.4
Order Index Order (5)	Method MAML-SEQ AGEM Replay OML-ER	$\begin{array}{c} \textbf{AGN} \\ A_{\rm CL,1} \\ 27.9 \\ 0.1 \\ 7.6 \\ 58.4 \end{array}$	$ ews F_1 56.1 89.9 82.4 27.7 $	Yel A _{CL,2} 20.4 55.8 56.5 43.4		Ama $A_{\rm CL,3}$ 19.456.157.243.9		Overall Forgetting 106.5 93.7 84.4 41.2
Order Index Order (5)	Method MAML-SEQ AGEM Replay OML-ER PMR	$\begin{array}{c} \textbf{AGN} \\ A_{\rm CL,1} \\ 27.9 \\ 0.1 \\ 7.6 \\ 58.4 \\ 0.0 \end{array}$	$ ews F_1 56.1 89.9 82.4 27.7 85.1 $	Yel A _{CL,2} 20.4 55.8 56.5 43.4 28.2		Ama: A _{CL,3} 19.4 56.1 57.2 43.9 31.4		Overall Forgetting 106.5 93.7 84.4 41.2 106.7
Order Index Order (5)	Method MAML-SEQ AGEM Replay OML-ER PMR Ours.(<i>all</i>)	AGN A _{CL,1} 27.9 0.1 7.6 58.4 0.0 68.0	$ ews F_1 56.1 89.9 82.4 27.7 85.1 14.3 \\ $	Yel A _{CL,2} 20.4 55.8 56.5 43.4 28.2 38.2	$\begin{array}{c} \mathbf{p} \\ F_2 \\ \hline 28.8 \\ 2.7 \\ 2.0 \\ 6.8 \\ 15.0 \\ 5.2 \end{array}$	Ama: A _{CL,3} 19.4 56.1 57.2 43.9 31.4 37.6		Overall Forgetting 106.5 93.7 84.4 41.2 106.7 21.5

Table 8: Per-task and Overall Forgetting on 3 Tasks

Table 9: Per-task and Overall Accuracy with Different Memory Constraints

Memory Constraint (B)	Method	Yelp	AGNews	Amazon	Overall Accuracy
	AGEM	57.3	0.0	58.0	38.5
B = 27	Replay	57.1	1.6	57.3	38.7
	OML-ER	49.8	19.4	49.1	39.4
	PMR	39.1	0.4	39.6	26.4
B = 27 (3 per class)	Ours.(all)	44.0	67.6	43.8	51.8
	Ours.(random)	36.2	63.6	35.1	44.9
	AGEM	57.0	1.7	56.7	38.5
B = 45	Replay	57.3	2.1	57.8	39.0
	OML-ER	48.0	13.1	49.3	36.8
	PMR	41.6	0.2	38.9	26.9
B = 45 (5 per class)	Ours.(all)	44.7	79.5	44.5	56.3
	Ours.(random)	41.3	74.3	40.7	52.1
	AGEM	57.2	0.0	57.9	38.4
B = 63	Replay	56.9	3.3	57.5	39.3
	OML-ER	52.7	24.6	50.3	42.5
	PMR	38.4	3.6	36.1	26.0
B = 63 (7 per class)	Ours.(all)	55.1	86.2	52.2	64.6
	Ours.(random)	40.9	69.3	39.3	49.8
	AGEM	54.9	4.7	55.4	38.4
	Replay	56.7	13.3	57.5	42.5
B = 34,500 (Unlimited)	OML-ER	46.8	49.2	47.3	47.8
D = 54,500 (Ommined)	PMR	-	_	-	—
	Ours.(all)	—	-	-	—
	Ours.(random)	55.1	86.2	52.2	64.6



Figure 1: The unigram distribution of the selected query samples. For simplicity, X-axis is the unigram index and Y-axis is the number of each unigram, where it is in the range [1, 26). Note that the size of selected sample set is fixed.