MoCoKGC: Momentum Contrast Entity Encoding for Knowledge Graph Completion

Anonymous ACL submission

Abstract

In recent years, a multitude of studies have 002 aimed to enhance the capabilities of pretrained language models (PLMs) for handling Knowledge Graph Completion (KGC) tasks by integrating structural information from knowledge graphs. However, existing approaches have not effectively combined the structural 007 attributes of knowledge graphs with textual descriptions of entities to generate entity encodings. To address this issue, this paper proposes MoCoKGC (Momentum Contrast Entity Encoding for Knowledge Graph Completion), in-012 corporating three primary encoders: the entityrelation encoder, the entity encoder, and the 015 momentum entity encoder. Through a slowly updated entity encoding mechanism, we main-017 tain a negative sample queue and characterize the entity neighborhood. On the standard evaluation metric, Mean Reciprocal Rank (MRR), the MoCoKGC model demonstrates superior performance, achieving a 7.1% improvement on the WN18RR dataset and an 11% improvement on the Wikidata5M dataset, while also surpassing the current best model on the FB15k-237 dataset. Through a series of experiments, this paper deeply studies the role and contribution of each component and parameter of the model.

1 Introduction

037

041

As an important method of knowledge representation, the fundamental building blocks of a knowledge graph are factual triples, such as (Steve Jobs, founded, Apple Inc.). However, the process of constructing knowledge graphs, whether manually or through automation, inevitably leads to the presence of many missing triplets within the knowledge graph. Therefore, the knowledge graph completion task (KGC) aims to complete the missing triples.

Among the numerous methods for KGC, knowledge graph embedding (KGE) techniques are the most classical and widely adopted. The core idea behind these methods is to generate embedding vectors for entities and relationships and use different scoring functions to predict the missing triplets (Bordes et al., 2013; Dettmers et al., 2018; Sun et al., 2019; Balazevic et al., 2019). Building upon this, some studies introduce the structure of knowledge graphs as supplementary information in the reasoning process to improve prediction accuracy (Schlichtkrull et al., 2018; Vashishth et al., 2020; Chen et al., 2021). 042

043

044

047

048

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

078

079

081

082

In recent years, researchers have begun exploring the integration of Pre-trained Language Models (PLMs) into the task of KGC, aiming to improve accuracy through the textual descriptions of entities and relationships. Models based on pre-trained language encoders can be broadly categorized into three types: (1) Cross-encoder models (Yao et al., 2019); (2) Bi-encoder models (Wang et al., 2021a, 2022);(3) Single-encoder models (Liu et al., 2022b; Chen et al., 2023a). Although cross-encoder models fully utilize the semantic information of triplets, their performance is often not as good as other methods due to the high cost of obtaining negative samples during training. Bi-encoder models, by separating the tail entity, not only preserve the textual information of the tail entity but also effectively increase the number of negative samples. Single-encoder models, despite removing the textual information of the tail entity, demonstrate unique advantages by acquiring a large number of negative samples for tail entities and introducing graph information through entity embeddings.

In order to preserve the textual information of entities while flexibly integrating entity encoding into the training and prediction phases of the model, this study adopts the momentum contrastive learning mechanism (He et al., 2020), thereby proposing the MoCoKGC model. This model draws inspiration from the Bi-encoder architecture, equipped with a head entity-relation encoder and an entity

encoder. The distinction lies in that, within the MoCoKGC model, the update of entity encoding 084 relies on a momentum entity encoder, rather than directly utilizing the entity encoder, thus achieving a smoothing of the entity encoding update process. Consequently, the MoCoKGC model exhibits two significant features compared to other methods: (1) Entity Queue. This queue is maintained in the order of entity encodings generated by the momentum entity encoder, providing the model with a rich source of negative tail entity samples; (2) Reutilization of entity encoding. Under this mechanism, entity encoding, as a part of integrating the structural information of the knowledge graph, is reimported into the encoder. Through this approach, MoCoKGC effectively addresses the challenges of negative sample generation and the integration of textual and structural information. 100

In terms of experimental validation, the Mo-CoKGC model not only successfully preserved the textual information of entities but also flexibly integrated entity encoding into the model's training and prediction processes. Its performance on standard datasets WN18RR, FB15k-237, and Wikidata5M demonstrates the model's superior capabilities: on the WN18RR dataset, in Mean Reciprocal Rank (MRR), the model achieved a 7.1% improvement, an 11% increase on the Wikidata5M dataset, and surpassed previous models on the FB15k-237 dataset. Furthermore, comparative experimental data reveal that MoCoKGC effectively overcomes the inconsistency issues exhibited by PLMs when dealing with sparse and dense knowledge graphs (Wang et al., 2022).

2 Related Work

101

102

103

104

105

106

108

109

110

111

112

113

114

115

116

117

Knowledge Graph Completion (KGC) tasks is to 118 predict missing triplet information within a knowl-119 edge graph. In the domain of knowledge graph 120 embeddings, a typical method is TransE (Bordes 121 et al., 2013), which utilizes the Euclidean distance 122 between the sum of head entity and relation em-123 beddings and the tail entity embedding as a scoring 124 function. The RotatE (Sun et al., 2019) method is 125 based on the core concept of interpreting relations 126 in the knowledge graph as rotational operations in 128 complex space. In tensor decomposition models such as Tucker (Balazevic et al., 2019), the knowl-129 edge graph is viewed as a three-dimensional tensor 130 and is realized by decomposing it into matrices for 131 entities, relations, and a core tensor. Graph-based 132

approaches, such as R-GCN (Schlichtkrull et al., 2018), address the problem of relation learning by introducing weight matrices for different types of relations in graph neural networks, thereby capturing the unique semantics of each relation type. CompGCN (Vashishth et al., 2020) encodes relations and replaces the relation transformation matrix by combining entity and relation embeddings.

133

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

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

Methods based on Pre-trained Language Models (PLMs), such as KG-BERT (Yao et al., 2019), concatenate the descriptions of head entities, relations, and tail entities, and directly obtain the triplet score by inputting it into the BERT (Devlin et al., 2019) model. StAR (Wang et al., 2021a) adopts a dual-encoder architecture, which significantly reduces the inference time overhead of language models. SimKGC (Wang et al., 2022) introduces a contrastive learning approach. Meanwhile, CSProm-KG (Chen et al., 2023a) enhances the model's understanding of the knowledge graph structure by using entity and relation embeddings as prompts for the language model.

Prompt Tuning has emerged as a strategy aimed at significantly improving the performance of PLMs by adding prompt tokens to the input. This approach was initially developed to address the challenge of fine-tuning Large Language Models (LLMs) for downstream tasks (Brown et al., 2020). Further developments such as Lester et al. (2021) train the model with special tokens as soft prompts, and p-tuning v2 (Liu et al., 2022a) introduces prompt parameters at every layer of the model, utilizing this deep prompting to enhance training effectiveness. Recently, in KGC tasks, researchers have improved the performance of PLMs through immediate learning, Lv et al. (2022) adding prompt templates and soft prompts to the input, while Chen et al. (2022) and Liu et al. (2022b) specified different soft prompt tokens for different types of relations. Chen et al. (2023a) not only introduced relation prompts but also added entity co-prompting to the model. This paper views relations as deep prompt parameters and introduces entity neighborhood prompts, achieving the objective of leveraging both textual descriptions and knowledge graph structural information.

Contrastive Learning by differentiating positive and negative sample features, learns distinctive feature representations and has been successfully

applied in multiple domains, including computer 185 MoCo (He et al., 2020) proposed a vision. 186 momentum-based contrastive learning method, effectively solving the problem of sample pair construction in unsupervised learning by building a dynamically changing encoder queue. SimCLR 190 (Chen et al., 2020), as a method of visual represen-191 tation learning, significantly improved the perfor-192 mance of image recognition tasks through large-193 scale unsupervised contrastive learning. In KGC 194 tasks, SimKGC (Wang et al., 2022) treats tail entities as positive and negative samples, implementing 196 efficient training through three different simple neg-197 ative sampling strategies. This paper combines the 198 momentum contrast method of MoCo, utilizing it 199 while dynamically updating entity encodings.

3 Methodology

3.1 Notations

201

203

207

208

211

212

213

214

215

216

217

218

219

220

222

227

229

232

Knowledge Graphs (KGs) represent a crucial data structure for organizing and storing factual relationships in the real world, which can be formally represented as $\mathcal{G} = \{\mathcal{E}, \mathcal{R}, \mathcal{T}\}$. Here, \mathcal{E} and \mathcal{R} denote the sets of entities and relationships, respectively. $\mathcal{T} = \{(h, r, t)\} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$ defines a set of triples, each comprising a head entity (h), a relation (r), and a tail entity (t). The task of KGC aims to fill in missing triples within a knowledge graph, with link prediction as its core task. This task focuses on predicting the missing entity part in given triples (h, r, ?) or (?, r, t). This task is typically accomplished by generating a ranked list of all possible candidate entities, where the ranking is based on the likelihood of a candidate entity being the missing head or tail entity.

3.2 Neighborhood Prompts

To integrate the structural information of knowledge graphs into pre-trained language models, this study proposes a method that utilizes the neighborhood information of entities as prompts. In knowledge graphs, the neighborhood of an entity is defined as the directly connected entities and their corresponding relations, formally represented as follows:

$$N(e) = \{ (e_i, r_i) | (e_i, r_i, e) \in \mathcal{T} \}$$
(1)

Given the variation in the neighborhood size of entities within knowledge graphs, this research introduces a parameter— σ —to standardize the dimension of sampled neighborhood information. Specifically, when the neighborhood size of an entity exceeds the set σ , a corresponding number of entity-relation pairs are randomly extracted from this neighborhood to meet the σ ; conversely, if the neighborhood size is smaller than the σ , specific padding tokens (pad tokens) are introduced to fill up to the σ . This treatment ensures the dimensional consistency of neighborhood information across all entities, facilitating subsequent processing. To obtain neighborhood prompts, entities and relations within the neighborhood are summed and then processed through a Multilayer Perceptron (MLP): 233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

268

269

270

271

272

273

274

275

276

277

278

279

280

 $\boldsymbol{p}_{N(e)} = \mathrm{MLP}([\boldsymbol{e}_0 + \boldsymbol{r}_0, ..., \boldsymbol{e}_{\boldsymbol{\sigma}} + \boldsymbol{r}_{\boldsymbol{\sigma}}]) \quad (2)$

3.3 Model Architecture

The MoCoKGC model is comprised of three primary components: the entity-relation encoder, the entity encoder, and the momentum entity encoder, as illustrated in Figure 1. The principal duties of these encoders are to generate encodings for the head entity and its relations, update the momentum entity encoder and pseudo-entity encodings, and produce entity encodings, respectively. It is noteworthy that the generation of entity encodings is dependent on the slower-updating momentum entity encoder, rather than the entity encoder . The following provides a detailed exposition of the functionality and operational flow of these three encoders.

Within the entity-relation encoder, the process initiates by aggregating the encodings of all entities within the vicinity of the head entity and their corresponding relations, followed by processing through a Multilayer Perceptron (MLP) to obtain neighborhood prompt information. Subsequently, the description of the head entity, the relation description, and the neighborhood prompt information are concatenated and inputted into a Transformer encoder. To more effectively amalgamate various types of information, we employ the p-tuning v2 strategy, as referenced in (Liu et al., 2022a), introducing the relation as deep prompt information at every layer of the Transformer encoder. The encoding of the head entity-relation is acquired through a pooling layer followed by normalization. This procedure can be formalized as:

$$\boldsymbol{h_r} = \text{ER_Encoder}(\boldsymbol{d}(h), \boldsymbol{d}(r), \boldsymbol{p}_{N(h)}, \boldsymbol{p}_r)$$
 (3)

It is important to highlight that the relation encoding within the neighborhood and the relation parameters in the deep prompts are not identical.



Figure 1: The MoCoKGC framework primarily consists of three encoders: the entity-relation encoder, the entity encoder, and the momentum entity encoder. As illustrated, the momentum entity encoder does not directly participate in the gradient backpropagation process; its parameter updates are based on the Exponential Moving Average (EMA) strategy. MoCoKGC updates entity encodings \mathcal{E} using a momentum entity encoder and augments the number of negative samples by maintaining an entity queue. Importantly, all entity encodings required for generating neighborhood prompts are sourced from \mathcal{E} .

In contrast to the entity-relation encoder, the entity encoder lacks the input of relation descriptions and relation cues, which can be represented as:

$$\bar{\boldsymbol{t}} = \text{E}_{\text{Encoder}}(\boldsymbol{d}(t), \boldsymbol{p}_{N(t)})$$
 (4)

The input to the momentum entity encoder is identical to that of the entity encoder.

$$\boldsymbol{t} = \text{ME}_\text{Encoder}(\boldsymbol{d}(t), \boldsymbol{p}_{N(t)})$$
(5)

In every iteration, newly generated entity encodings are utilized to update the old entity encodings. Uniquely, the parameters of this encoder are not updated through backpropagation but evolve iteratively based on the parameter updates of the entity encoder after each iteration. The iterative formula is as follows:

287

290

294

298

$$\theta_{ME} = m\theta_{ME} + (1-m)\theta_E \tag{6}$$

Where θ_{ME} and θ_E denote the parameters of the momentum entity encoder and the entity encoder, respectively. $m \in [0, 1]$ represents the momentum

coefficient. This process is also referred to as the Exponential Moving Average (EMA). The larger the momentum coefficient m, the slower the update of θ_{ME} .

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

It is crucial to highlight that the neighborhood representation of entity encodings, employed by all three encoders, is shared and generated by the momentum entity encoder. Conversely, the relation encoding is unique to each model component and is not shared.

3.4 Negative Sampling

In-batch Negatives like most contrastive learning methods, our study uses tail entities from within the same batch as negative samples. In our approach, we not only utilize entity encodings as positive and negative examples but also integrate them into the representations of their respective neighborhoods. Therefore, we not only process entities present in the given batch of triples but also sequentially generate additional entity encodings and pseudo-entity encodings.

375

376

377

378

379

380

381

383

384

385

386

387

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

368

$$scores = \{ \boldsymbol{h}_{\boldsymbol{r}} \boldsymbol{t}_{\boldsymbol{i}}^T | t_{\boldsymbol{i}} \in \mathcal{E} \}$$
 (9)

In terms of time complexity, the time complexity of this study in the test set is consistent with that of most KGC models, which is $|\mathcal{T}_{test}|$.

processes, a final update of all entities will be con-

For the prediction inference of KGC, it is only

necessary to generate the head entity-relation en-

coding through the entity-relation encoder, and

then multiply it with all entity encodings to obtain

the predictive scores for all entities.

4 Experiments

ducted.

Dataset	# entity	relation	# train	# valid	# test
WN18RR	40,943	11	86,835	3,034	3,134
FB15k-237	14,541	237	272,115	17,535	20,466
Wikidata5M	4,594,485	822	20,614,279	5,163	5,163

Table 1: Summary statistics of benchmark datasets.

4.1 Experimental Setup

Dataset Evaluation In this study, three benchmark datasets were utilized to assess the performance of the proposed model, specifically: WN18RR, FB15k-237, and Wikidata5M. Table 1 presents the detailed distribution of these datasets. The WN18RR dataset (Dettmers et al., 2018) is constructed based on the WordNet knowledge base (Miller, 1998), aimed at link prediction tasks, containing entities represented by English phrases and their semantic relationships. The FB15k-237 dataset (Toutanova et al., 2015) is a subset derived from the Freebase knowledge base (Bollacker et al., 2008), encompassing entities in the real world and their interrelations. The Wikidata5M dataset (Wang et al., 2021b) is a large-scale knowledge graph dataset, integrating information from the Wikidata knowledge graph and Wikipedia pages, providing Wikipedia page descriptions for each entity. Compared to WN18RR and FB15k-237, the Wikidata5M dataset surpasses them by two orders of magnitude in both the number of entities and triples, indicating its larger scale and complexity.

Evaluation Metrics In the task of KGC, the assessment of model performance is primarily achieved by measuring the ranking of target triples among all potential triples' scores. This study adopts the commonly used evaluation metrics in previous research, including Hits@1, Hits@3,

Entity Queue is maintained by MoCoKGC throughout the training process to generate a larger pool of negative sample entities. Unlike conventional queues, the entity elements in this queue are unique. If an entity that is about to be enqueued is already present in the queue, it is first dequeued and then enqueued again. This mechanism ensures that a greater variety of different entities can be stored while the queue length remains fixed.

3.5 Training and Inference

321

322

326

327

334 335

339

340

341

347

351

352

357

361

364

367

During the training phase of the model, considering that the neighborhood sampling of the head entity may contain the tail entity, and similarly, the neighborhood sampling of the tail entity may include the head entity, this study adopts a target link dropout strategy after the neighborhood sampling process. Specifically, the target links existing in the neighborhoods of the head and tail entities are discarded and replaced with a padding token. This measure aims to ensure that training target data is not leaked into the model, thereby affecting the model's generalization capability. In addition to the dropout of target links, to enhance the diversity of neighborhood information, this study also introduces a mechanism to randomly drop entity-relation pairs in the neighborhood with a certain probability.

As illustrated in Figure 1, the process of loss calculation in this study involves multiplying the generated head entity relation encoding h_r with both the pseudo-entity encoding \bar{t} and the actual entity encoding t, based on which the loss is calculated. This study employs the same method of calculating the loss function as SimKGC (Wang et al., 2022), use InfoNCE loss with additive margin (Chen et al., 2020; Yang et al., 2019):

$$\mathcal{L}(\boldsymbol{h_r}, \boldsymbol{t}) = -\log \frac{e^{(\boldsymbol{h_r} \boldsymbol{t}^T - \boldsymbol{\gamma})/\tau}}{e^{(\boldsymbol{h_r} \boldsymbol{t}^T - \boldsymbol{\gamma})/\tau} + \sum_{i=1}^{|\mathcal{N}|} e^{(\boldsymbol{h_r} \boldsymbol{t'_i}^T)/\tau}}$$
(7)

Where γ is the margin coefficient greater than 0, $\tau \in [0, 1]$ is the temperature coefficient and \mathcal{N} is all negative sample entities. Based on the formula 7, the final loss function can be expressed as:

$$loss = \mathcal{L}(\boldsymbol{h}_{\boldsymbol{r}}, \boldsymbol{t}) + \mathcal{L}(\boldsymbol{h}_{\boldsymbol{r}}, \bar{\boldsymbol{t}})$$
(8)

To enhance the update frequency of entity encodings, this study not only updates a portion of entity encodings at each iteration but also separately utilizes the momentum entity encoder for inference after a certain number of iterations, to achieve updates of all entities. After completing all training

	WN18RR		FB15k-237			Wikidata5M						
	MRR	Hits@1	Hits@3	Hits@10	MRR	Hits@1	Hits@3	Hits@10	MRR	Hits@1	Hits@3	Hits@10
Knowledge graph embedding method												
TransE (Bordes et al., 2013) [◊]	0.243	0.043	0.441	0.532	0.279	0.198	0.376	0.441	0.253	0.170	0.311	0.392
DistMult (Yang et al., 2015) [◊]	0.444	0.412	0.470	0.504	0.281	0.199	0.301	0.446	0.253	0.209	0.278	0.334
ComplEx (Trouillon et al., 2016) [◊]	0.449	0.409	0.469	0.530	0.278	0.194	0.297	0.450	<u>0.308</u>	<u>0.255</u>	-	<u>0.398</u>
R-GCN (Schlichtkrull et al., 2018) [†]	0.123	0.080	0.137	0.207	0.164	0.100	0.181	0.300	-	-	-	-
ConvE (Dettmers et al., 2018) [†]	0.456	0.419	0.470	0.531	0.312	0.225	0.341	0.497	-	-	-	-
RotatE (Sun et al., 2019) [◊]	0.476	0.428	0.492	0.571	0.338	0.241	0.375	0.533	0.290	0.234	0.322	0.390
TuckER (Balazevic et al., 2019)	0.470	0.443	0.482	0.526	0.358	0.266	0.394	0.544	-	-	-	-
CompGCN (Vashishth et al., 2020)	0.479	0.443	0.494	0.546	0.355	0.264	0.390	0.535	-	-	-	-
HittER (Chen et al., 2021)	<u>0.503</u>	0.462	0.516	<u>0.584</u>	0.373	0.279	0.409	0.558	-	-	-	-
N-Former (Liu et al., 2022b)	0.486	0.443	0.501	0.578	0.372	0.277	<u>0.412</u>	0.556	-	-	-	-
PLM-Based method												
KG-BERT (Yao et al., 2019)	0.216	0.041	0.302	0.524	-	-	-	0.420	-	-	-	-
StAR (Wang et al., 2021a)	0.401	0.243	0.491	0.709	0.296	0.205	0.322	0.482	-	-	-	-
KEPLER(Wang et al., 2021b) [◊]	-	-	-	-	-	-	-	-	0.210	0.173	0.224	0.277
KG-S2S (Chen et al., 2022)	0.574	0.531	0.595	0.661	0.336	0.257	0.373	0.498	-	-	-	-
N-BERT (Liu et al., 2022b)	0.583	0.529	0.607	0.686	0.381	0.287	0.420	0.562	-	-	-	-
SimKGC (Wang et al., 2022)	0.671	0.585	0.731	0.817	0.333	0.246	0.362	0.510	0.358	0.313	0.376	0.441
CSProm-KG (Chen et al., 2023b)	0.575	0.522	0.596	0.678	0.358	0.269	0.393	0.538	<u>0.380</u>	<u>0.343</u>	<u>0.399</u>	<u>0.446</u>
GHN (Qiao et al., 2023)	<u>0.678</u>	<u>0.596</u>	<u>0.719</u>	0.821	0.339	0.251	0.364	0.518	0.364	0.317	0.380	0.453
Ensemble method												
StAR(Self-Adp) (Wang et al., 2021a)	0.551	0.459	0.594	0.732	0.365	0.266	0.404	0.562	-	-	-	-
CoLE (Liu et al., 2022b)	<u>0.587</u>	<u>0.532</u>	<u>0.608</u>	<u>0.694</u>	<u>0.389</u>	<u>0.294</u>	<u>0.430</u>	<u>0.574</u>	-	-	-	-
MoCoKGC(Ours)	0.742	0.665	0.792	0.881	0.391	0.296	0.431	0.580	0.490	0.435	0.517	0.591

Table 2: Experimental results for various baseline methods on WN18RR, FB15k-237, and Wikidata5M datasets. [†]: Results are sourced from Wang et al. (2021a); \diamond : Results are sourced from Chen et al. (2023b). The best methods are highlighted in bold, with the most effective methods in each category underscored for emphasis. The results for the MoCoKGC model are reported as the average of three experimental runs.

Hits@10, and MRR. The Hits@k metric measures the frequency with which the target triple appears among the top k triples with the highest scores, while the MRR is the average of the reciprocal ranks of the target triples. To enhance the accuracy and fairness of the evaluation, we employed the filtered ranking setting proposed by (Bordes et al., 2013), which eliminates potential ranking biases by excluding all possible triples (h, r, ?)or $(t, r^{-1}, ?)$ that already exist in the training set. Furthermore, following the random evaluation protocol suggested by Sun et al. (2020), we accurately assess model performance. The final metric is the average of the forward prediction (h, r, ?) and the backward prediction $(t, r^{-1}, ?)$.

Implementation Details To ensure the comparability of the results of this study with existing research, we selected the "bert-base-uncased" version of the BERT model as the Transformer encoder for this research. The experimental environment was configured on a GeForce RTX 4090 GPU server equipped with 24GB of graphics card memory. utilizing the AdamW optimizer for model training. The learning rate was set to 5×10^{-5} . The batch size was selected from the set {256, 512, 1024}. The range of the momentum coefficient *m* was chosen from {0, 0.5, 0.9, 0.99, 0.999}. The neighborhood sampling size σ was selected from

the set {256, 512, 1024}. The length of the maintained entity queue was chosen from the set {512, 1024, 2048, 4096, 8192, 16384, 32768}. For the WN18RR, FB15k-237, and Wikidata5M datasets, training was conducted for 30, 3, and 1 epochs, respectively. For further details, please refer to Appendix A.



Figure 2: MRR performance of different models on WN18RR and FB15k-237 datasets.

4.2 Main Results

On the WN18RR, FB15k-237, and Wikidata5M datasets, we compared the MoCoKGC model with other leading models, as shown in Table 2. The experimental results demonstrate that MoCoKGC

achieved state-of-the-art performance across all evaluation metrics. Notably, on the WN18RR and Wikidata5M datasets, MoCoKGC realized significant improvements of 7.1% (from 0.671 to 0.742) and 11% (from 0.343 to 0.399), respectively. As a method based on pre-trained language models (PLM-Based), MoCoKGC also achieved a 1.0% performance improvement (from 0.381 to 0.391) on the FB15k-237 dataset, surpassing the previous best ensemble learning approach.

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

Furthermore, we conducted a separate analysis on the MRR values of models that performed well on the WN18RR and FB15k-237 datasets, as depicted in Figure 2. The analysis revealed that knowledge graph embedding methods exhibited relatively balanced performance on these two datasets (i.e., models that performed well on WN18RR also excelled on FB15k-237). In PLM-based models, SimKGC and GHN exhibit significant performance improvements on the WN18RR dataset, yet they lag on the FB15k-237 dataset. We attribute this phenomenon to SimKGC's use of entity descriptions, generating entity encodings through an entity encoder, and the absence of knowledge graph structural information during inference. MoCoKGC successfully addressed the inconsistency in performance of PLM-based models on these two datasets.

> On the larger Wikidata5M dataset, the performance improvement of MoCoKGC was especially pronounced, which is closely related to the rich entity textual descriptions and significant knowledge graph structure within the Wikidata5M dataset. Our proposed MoCoKGC model, as a PLM-based method, not only integrates the entity encoder from SimKGC (Wang et al., 2022) but also, like models such as CoLE (Liu et al., 2022b) and CSProm-KG (Chen et al., 2023b), incorporates the structure of knowledge graphs (e.g., relation and neighborhood prompts) into the model. This effectively combines the advantages of textual descriptions with the knowledge graph structure.

4.3 Ablation Studies

Model	MRR	Hits@1	Hits@10
MoCoKGC w/o neighborhood prompt MoCoKGC w/o relation prompt	0.696 0.597	0.614 0.476	0.845 0.818
MoCoKGC	0.742	0.665	0.881

Table 3: Ablation Study regarding important components in MoCoKGC on the benchmark of WN18RR.

Model	MRR	H@1	H@10
MoCoKGC w/o neighborhood prompt MoCoKGC w/o relation prompt	0.385 0.327	0.292 0.242	0.570 0.496
MoCoKGC	0.391	0.296	0.580

Table 4: Ablation Study regarding important components in MoCoKGC on the benchmark of FB15k-237.

Structural component. In this study, we conducted an in-depth analysis of the impact on performance by removing two key components from the MoCoKGC model: the neighborhood prompt and the relation prompt. Specific experiments were carried out on two benchmark knowledge graph datasets, WN18RR and FB15k-237, with their results presented in Table 3 and Table 4, respectively. The experimental outcomes reveal that the absence of neighborhood prompts diminished the model's performance on the WN18RR and FB15k-237 datasets by 4.9% (from 0.745 to 0.696) and 0.6% (from 0.391 to 0.385), respectively. This strongly evidences that the model indeed learns relevant information about neighborhood structures from the neighborhood prompts. More notably, the removal of relation prompts led to a substantial performance decline of 14.5% (from 0.745 to 0.597) and 6.4% (from 0.391 to 0.327) on the WN18RR and FB15k-237 datasets, respectively. This significant performance drop was even more pronounced than that of the SimKGC model.

This phenomenon suggests that the simple reuse of entity encodings might interfere with the encoder's effective capture of deep semantic information about entities and their relations. To overcome this issue, the introduction of relation prompts is crucial for restoring and enhancing the synergistic effect of textual semantics and knowledge graph structural information within PLMs.

m	0	0.5	0.9	0.99	0.999
MRR (WN18RR)	0.727	0.728	0.728	0.733	0.742
MRR (FB15k-237)	0.369	0.367	0.375	0.380	0.391

Table 5: Demonstrates the MRR of MoCoKGC on the datasets WN18RR and FB15k-237, with varying momentum coefficient m used during training.

Momentum coefficient. In Table 5, we present the results of the MRR for models trained with different momentum coefficients m on the

525

526

527

493

WN18RR and FB15k-237 datasets. Analysis indicates that higher momentum coefficients mcan stably enhance model performance, whereas lower momentum coefficients m have not shown significant improvement in performance. This experimental outcome aligns with our initial rationale for employing a momentum entity encoder, which is to introduce a steady yet gradual entity encoding update mechanism, in the hope of achieving performance improvement.

528

534

537

538

539

540

546

547

550



Figure 3: Variation of MRR with entity queue size on WN18RR in MoCoKGC



Figure 4: Variation of MRR with entity queue size on FB15k-237 in MoCoKGC

Entity queue size. In the training framework of MoCoKGC, a pivotal component is the maintenance of a dynamic entity queue, aimed at accumulating and leveraging a broader spectrum of negative tail entity samples throughout the training process. To investigate the impact of the entity queue, we examined how variations in the size of the entity queue influence model performance. As illustrated in Figures 3 and 4, the Mean Reciprocal Rank (MRR) on WN18RR and FB15k-237 varies with different entity queue sizes. The results demonstrate a consistent upward trend in the MRR metric as the size of the entity queue increases. This indicates that expanding the entity queue significantly augments the quantity of effective negative entity samples, thereby exerting a positive impact on model performance.



Figure 5: Variation of MRR with training set size on Wikidata5M in MoCoKGC, with comparative final MRR sesults from SimKGC and CSProm-KG on the entire training set.

Impact of Training Set Size. During the training process of the MoCoKGC model, we observed that the model could achieve commendable performance even with a limited amount of training data. As depicted in Figure 5, for comparison purposes, we presented the training outcomes of the SimKGC and CSProm-KG models on the complete training set in dashed lines. Notably, when utilizing only 20% of the training data, the MRR of the MoCoKGC model could reach 0.460. This result significantly surpasses the final performance of the other two methods. This finding underscores the exceptional generalization capability of MoCoKGC in scenarios of data scarcity.

5 Conclusion

This study proposes MoCoKGC, a novel KGC model that leverages momentum contrastive learning in conjunction with PMLs. By expanding the pool of negative samples, it further enhances KGC through the aggregation of entity textual descriptions and their structural information. The Mo-CoKGC model demonstrated superior performance across multiple datasets. Furthermore, we further validated the critical role of its constituent components and parameter configurations. Future work will focus on adapting MoCoKGC for open knowledge graphs to better manage the emergence of new entities.

583

6 Limitations

585

587

588

589

591

597

607

608

611

612

613

614

615

616

617

618

619

625

626

627

628

629

630

631

634

The MoCoKGC model relies on pre-trained language models to integrate textual representations with the structure of knowledge graphs. This results in an increase in training time and memory consumption as the length of the structure input into the model increases. In response, MoCoKGC opts for a compromise by sampling the neighborhoods of entities, rather than aggregating the entire knowledge graph structure as done by R-GCN (Schlichtkrull et al., 2018) and CompGCN (Vashishth et al., 2020). Moreover, the random sampling does not take into account the varying importance of different links within the neighborhood. This leads to the model predictions being more focused on the features within the sampled neighborhoods. In the future, we plan to address this issue.

References

- Ivana Balazevic, Carl Allen, and Timothy M. Hospedales. 2019. Tucker: Tensor factorization for knowledge graph completion. In <u>Proceedings</u> of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019, pages 5184–5193. Association for Computational Linguistics.
- Kurt D. Bollacker, Colin Evans, Praveen K. Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In Proceedings of the ACM <u>SIGMOD International Conference on Management</u> of Data, SIGMOD 2008, Vancouver, BC, Canada, June 10-12, 2008, pages 1247–1250. ACM.
- Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko.
 2013. Translating embeddings for modeling multirelational data. In <u>Advances in Neural Information</u> <u>Processing Systems 26: 27th Annual Conference</u> <u>on Neural Information Processing Systems 2013.</u> <u>Proceedings of a meeting held December 5-8, 2013,</u> <u>Lake Tahoe, Nevada, United States, pages 2787– 2795.</u>
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei.

2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems
33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December
6-12, 2020, virtual.

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

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

- Chen Chen, Yufei Wang, Bing Li, and Kwok-Yan Lam. 2022. Knowledge is flat: A seq2seq generative framework for various knowledge graph completion. In Proceedings of the 29th International Conference on Computational Linguistics, COLING 2022, Gyeongju, Republic of Korea, October 12-17, 2022, pages 4005–4017. International Committee on Computational Linguistics.
- Chen Chen, Yufei Wang, Aixin Sun, Bing Li, and Kwok-Yan Lam. 2023a. Dipping plms sauce: Bridging structure and text for effective knowledge graph completion via conditional soft prompting. In Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023, pages 11489–11503. Association for Computational Linguistics.
- Chen Chen, Yufei Wang, Aixin Sun, Bing Li, and Kwok-Yan Lam. 2023b. Dipping plms sauce: Bridging structure and text for effective knowledge graph completion via conditional soft prompting. In <u>Findings</u> of the Association for Computational Linguistics: <u>ACL 2023, Toronto, Canada, July 9-14, 2023, pages</u> 11489–11503. Association for Computational Linguistics.
- Sanxing Chen, Xiaodong Liu, Jianfeng Gao, Jian Jiao, Ruofei Zhang, and Yangfeng Ji. 2021. Hitter: Hierarchical transformers for knowledge graph embeddings. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021, pages 10395– 10407. Association for Computational Linguistics.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. 2020. A simple framework for contrastive learning of visual representations. In Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event, volume 119 of Proceedings of Machine Learning Research, pages 1597–1607. PMLR.
- Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. 2018. Convolutional 2d knowledge graph embeddings. In Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018, pages 1811–1818. AAAI Press.
- 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

- 710
- 711 712
- 715 716 717
- 724
- 726
- 727 728 729
- 731
- 733 734 735 736

737 740

- 741 742
- 743
- 745 746

744

747

- 749
- 750
- 751

754

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.

- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020, pages 9726-9735. Computer Vision Foundation / IEEE.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 3045-3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022a. Ptuning: Prompt tuning can be comparable to finetuning across scales and tasks. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 61-68, Dublin, Ireland. Association for Computational Linguistics.
- Yang Liu, Zegun Sun, Guangyao Li, and Wei Hu. 2022b. I know what you do not know: Knowledge graph embedding via co-distillation learning. In Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17-21, 2022, pages 1329–1338. ACM.
- Xin Lv, Yankai Lin, Yixin Cao, Lei Hou, Juanzi Li, Zhiyuan Liu, Peng Li, and Jie Zhou. 2022. Do pre-trained models benefit knowledge graph completion? A reliable evaluation and a reasonable approach. In Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022, pages 3570-3581. Association for Computational Linguistics.
- George A Miller. 1998. WordNet: An electronic lexical database. MIT press.
- Zile Qiao, Wei Ye, Dingyao Yu, Tong Mo, Weiping Li, and Shikun Zhang. 2023. Improving knowledge graph completion with generative hard negative mining. In Findings of the Association for Computational Linguistics: ACL 2023, pages 5866-5878, Toronto, Canada. Association for Computational Linguistics.
- Michael Sejr Schlichtkrull, Thomas N. Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and Max Welling. 2018. Modeling relational data with graph convolutional networks. In The Semantic Web - 15th International Conference, ESWC 2018, Heraklion,

Crete, Greece, June 3-7, 2018, Proceedings, volume 10843 of Lecture Notes in Computer Science, pages 593-607. Springer.

755

756

757

758

759

760

761

762

763

764

765

768

771

772

776

778

779

780

781

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

- Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. 2019. Rotate: Knowledge graph embedding by relational rotation in complex space. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net.
- Zhiqing Sun, Shikhar Vashishth, Soumya Sanyal, Partha P. Talukdar, and Yiming Yang. 2020. A re-evaluation of knowledge graph completion methods. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020, pages 5516-5522. Association for Computational Linguistics.
- Kristina Toutanova, Danqi Chen, Patrick Pantel, Hoifung Poon, Pallavi Choudhury, and Michael Gamon. 2015. Representing text for joint embedding of text and knowledge bases. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015, pages 1499-1509. The Association for Computational Linguistics.
- Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. 2016. Complex embeddings for simple link prediction. 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 2071-2080. JMLR.org.
- Shikhar Vashishth, Soumya Sanyal, Vikram Nitin, and Partha P. Talukdar. 2020. Compositionbased multi-relational graph convolutional networks. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenRe-Learning view.net.
- Bo Wang, Tao Shen, Guodong Long, Tianyi Zhou, Ying Wang, and Yi Chang. 2021a. Structureaugmented text representation learning for efficient knowledge graph completion. In WWW '21: The Web Conference 2021, Virtual Event / Ljubljana, Slovenia, April 19-23, 2021, pages 1737-1748. ACM / IW3C2.
- Liang Wang, Wei Zhao, Zhuoyu Wei, and Jingming Liu. 2022. Simkgc: Simple contrastive knowledge graph completion with pre-trained language models. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022, pages 4281-4294. Association for Computational Linguistics.
- Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhengyan Zhang, Zhiyuan Liu, Juanzi Li, and Jian Tang. 2021b.

KEPLER: A unified model for knowledge embedding and pre-trained language representation. <u>Trans.</u> <u>Assoc. Comput. Linguistics</u>, 9:176–194.

812

813

814

815 816

817

818

823

827

829

830

831

832

833

835

837

839

840

843

- Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. 2015. Embedding entities and relations for learning and inference in knowledge bases. In <u>3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings.</u>
- Yinfei Yang, Gustavo Hernández Ábrego, Steve Yuan, Mandy Guo, Qinlan Shen, Daniel Cer, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil.
 2019. Improving multilingual sentence embedding using bi-directional dual encoder with additive margin softmax. In <u>Proceedings of the</u> <u>Twenty-Eighth International Joint Conference on</u> <u>Artificial Intelligence, IJCAI 2019, Macao, China,</u> <u>August 10-16, 2019, pages 5370–5378. ijcai.org.</u>
 - Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. KG-BERT: BERT for knowledge graph completion. CoRR, abs/1909.03193.

A Details on Implementation

# of GPUs	1
learning rate	5×10^{-5}
initial temperature τ	0.05
gradient clip	10
warmup steps	400
dropout	0.1
neighborhood dropout	0.1
weight decay	10^{-4}
InfoNCE margin	0.02
momentum coefficient m	0.999
pooling	mean

Table 6: Shared hyperparameters for MoCoKGC.

	WN18RR	FB15k-237	Wikidata5M
batch size	1024	256	1024
additional entity size	512	1024	512
max # of word tokens	64	128	64
neighborhood sampling size σ	16	128	32
entity queue size	16384	14541	16384
epochs	30	3	1

Table 7: Hyperparameters of the MoCoKGC model that are not shared across different datasets.

In this study, the hyperparameter settings were primarily aligned with the configuration strategy of SimKGC (Wang et al., 2022). As demonstrated in Table 6, we have listed the hyperparameter settings shared across all datasets. Concurrently, Table 7 showcases the specific hyperparameter configurations for the MoCoKGC model across different datasets.

Given that the experiments were conducted using a single GeForce RTX 4090 graphics card, and faced with memory capacity limitations, we employed gradient accumulation techniques to enable larger batch sizes. It is noteworthy that, due to the infeasibility of directly applying conventional gradient accumulation methods in the contrastive learning process, we first generate all necessary contrastive encodings for the three encoders using smaller batch sizes and disabling gradient saving during each accumulation step. Subsequently, we update the entity-relation encoder and the entity encoder using gradient accumulation techniques. To eliminate the potential randomness introduced by dropout operations, a random number is generated and recorded as the random seed during each gradient accumulation, and this seed is set every time an encoder is invoked. In addition to gradient accumulation, in the experiments on Wikidata5M, we stored the entity encodings in CPU memory rather than in GPU memory to reduce the usage of GPU memory.

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

881

882

883

During each training epoch, the MoCoKGC model runs on a single GeForce RTX 4090 graphics card, utilizing a configuration that includes four workers for data loading. The runtime varies depending on the dataset: it takes approximately 7 minutes for the WN18RR dataset, about 40 minutes for the FB15k-237 dataset, and roughly 65 hours for the Wikidata5M dataset.

Furthermore, drawing from the practices of SimKGC, we made the following adjustments to the textual descriptions of entities: (1) the names of neighboring entities in the training set are concatenated to the description of the entity, and the correct entities are dynamically excluded from the input text during the training process; (2) the descriptions of inverse relations are formed by appending the term "inverse" to the beginning of the original relation descriptions.

Our implementation is based on open-source project *transformers*¹.

¹https://github.com/huggingface/transformers