Learning Semantic Matching via Augment-normalized Language Model with Commix Dimensional Attention

Anonymous NAACL submission

Abstract

Semantic matching is a fundamental task in Natural Language Processing (NLP), which is widely used in information retrieval, recommendation, and other applications. Transformer-based pre-trained language models have achieved remarkable improvements in semantic matching. However, the transformer uses only one attention mechanism, which might not be optimal for semantic matching that relies on the modeling of complex relationships. In this paper, we propose the Commix Dimensional Attention(CDA) framework to enhance the ability of language models to capture the relationships between sen-015 tence pairs from diverse aspects by exploiting and commixing four complementary attention mechanisms. Building based upon the transformer architecture, the method adopts diverse types of attention functions to capture manifold types of interactive information and effectively fuses them with a well-designed selfinteractive augmentation layer and a normalized aggregation layer. Specifically, the CDA language model includes three key modules, 1) a commix dimensional attention module, 2) a self-interactive augmentation module, and 3) a normalized aggregation module. We apply the proposed CDA language model to conduct extensive experiments. Results show that the proposed model achieves consistent improvement on 10 well-studied semantic matching datasets.

1 Introduction

011

014

017

Semantic Sentence Matching (SSM) plays an important role in Natural Language Processing (NLP). SSM aims to compare two sentences and identify their semantic relationship. In recent years, with the development of pre-trained language models(PLMs), PLMs with attention are regarded as the core structure, such as Bert (Devlin et al., 2018), RoBERTa (Liu et al., 2019b). The PLMs gen-041 erally adopt large-scale training corpus and selfsupervised learning objectives to learn sentence



Figure 1: Examples of other PLMs like BERT, RoBERTa, or even chatGPT can not distinguish the semantics of similar texts well.

representation better. With the powerful context representation ability, they have achieved state-ofthe-art performance in semantic matching tasks. Recent work shows that using external knowledge to enhance the attention mechanism can further improve the performance of the model(Han et al., 2021). For example, SyntaxBERT(Bai et al., 2021) proposes to exploit the text sentence structure and enhance the model by adding syntactic information to the attention. By introducing synonym information to enhance the attention mechanism in the pre-trained language model, UERBERT(Xia et al., 2021) achieves significant performance improvements. These works show the importance of incorporating some inductive bias into the attention mechanism for text sequence learning. Besides, Large language models (LLMs) have revolutionized natural language task solving through prompting(Brown et al., 2020) and have demonstrated impressive capabilities in a variety of natural language

063

processing tasks. Large language models such as GPT3(Brown et al., 2020), GLM2(Zeng et al., 2022), chatGPT, LLaMa(Touvron et al., 2023), etc. have strong capabilities in the field of generative domain. However, with their huge model parameters and complex model structures, they also have strong abilities in semantic matching.

071

077

078

084

090

096

100

101

102

103

104

106

107

108

110

111

112

113

114

Most of the current language models are based on transformer architecture. The calculation of the attention score in the transformer is merely based on the dot-product to model the relationship between sentence pairs, which may not be optimal for the transformer-based language model. However, These models do not perform very well in distinguishing sentence pairs with high literal similarities. Figure 1 demonstrates a few cases suffering from this problem. Although the sentence pairs in this figure are semantically different, they are too similar in literally for those language models like BERT, RoBERTa, or even chatGPT to distinguish accurately. In addition, from case 3 and case 4 in this figure, there is a phenomenon of instability in the task of determining semantic similarity in chatGPT, and two opposite results are given for the same input sentence pair. At the same time, De-attention (Tay et al., 2019) and MwAN (Tan et al., 2018) work in the non-pre-trained model have verified the effectiveness of a flexible attention mechanism. But it is still not sure whether a more flexible attention model in the large-scale pretrained model works well on semantic matching and how to design an effective flexible attention model to enhance semantic matching.

To this end, in this paper, we propose a new model, named Commix Dimensional Attention Language Model (CDA-LM), Bert is a representative model of bidirectional language models, so we apply our CDA mechanism in Bert, as the attention in the first layer transformer of Bert is broad and uninformed (Xia et al., 2021), meanwhile, so as to avoid adding unnecessary parameters, The proposed model reforms the multi-head attention module of the first layer transformer of BERT by incorporating four complementary attention mechanisms, interactively augment and adaptively integrate four kinds of attention information for sentence matching. due to the fact that the attention information is captured from commix perspectives or views, we call our attention mechanism Commix Dimensional Attention. Specifically, it includes three modules:

1) Commix Dimensional Attention module. We analyze that the four dimensions of attention are complementary clues for sentence matching, which can capture different levels of information in the text sequence. We propose a commix dimensional attention model by considering four different attention mechanisms, including dot-product attention, additive attention, minus attention, and bilinear attention. Compared with the single attentionbased model, our framework can model the relationship between sentences from different dimensions through commix dimensional attention, so as to obtain more fine-grained matching information. 115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

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

2) Self-interactive Augmentation module. We observe that The representations obtained through different dimensional attention mechanisms do not interact well with input information, We propose to apply the self-interactive augmentation module to interact with the matching information together with each word in the sentence in each attention function. So in order to augment the interaction with input information and thus obtain better semantic representations.

3) Normalized Aggregation module. We find that the simple aggregation with fixed or average importance weights may destroy the learned knowledge of the pre-trained language model. We propose to adaptively aggregate the representations obtained by the self-interactive augmentation module, normalized aggregation combines the matching information of all attention functions. We apply the normalized aggregation module to aggregate the four representations adaptively.

In order to verify the effectiveness of our proposed model, we conducted intensive experiments on 10 datasets, including GLUE datasets such as QQP, MRPC, and SNLI, and fully studied datasets such as Sci and Twi. The results show that compared with BERT-base, CDA-LM achieves an absolute improvement of more than 2.2% avg, and is superior to other Bert-based models and large language models(LLMs) in more advanced technology and external data use.

The main contributions of this work can be summarized as follows:

- We provide an in-depth analysis of the feasibility of improving the attention mechanism in the pre-trained model and propose a new Commix Dimensional Attention Language Model (CDA-LM).
- The proposed CDA-LM effectively augments

166and aggregates four complementary attention167models, such that the intrinsic complex rela-168tionship between sentence pairs can be fully169discovered for effective semantic matching.

• Extensive experiments are conducted on 10 semantic matching datasets. The results show that the proposed CDA-LM achieves remarkable performance gain compared with BERT (with 2.2% improvements on average) and also outperforms the state-of-the-art external knowledge enhancement-based methods.

2 Related Work

170

171

172

173

174

175

176

177

178

179

180

182

183

184

186

187

188

190

191

192

195

196

199

203

207

208

211

212

213

214

215

Semantic Sentence Matching plays an important role in many applications, such as information retrieval (IR) and natural language inference (NLI).

Recently, the shift from neural network architecture engineering to large-scale pre-training has significantly improved NLP tasks, demonstrating the power of unsupervised pre-training. large-scale pretrained language models (PLMs) have boosted the performance of text semantic matching by making full use of massive text resources. Most of them are composed of multiple transformer layers(Vaswani et al., 2017) with multi-head attention and are pretrained with well-designed self-supervised learning objectives. Outstanding examples include Embedding from Language Models(ELMo) (Peters et al., 2018), Generative Pre-trained Transformers (GPT) (Radford et al., 2018), Bidirectional Encoder Representations from Transformers (BERT) (Devlin et al., 2018), and Generalized Auto-regressive Pre-training (XLNet) (Yang et al., 2019). Providing fine-grained contextual word embedding, these pre-trained models can be either easily applied to downstream tasks as encoders or directly finetuned for downstream tasks. As the most prominent model in recent years, BERT and many of its variants, including AlBERT (Lan et al., 2019), RoBERTa (Liu et al., 2019b), ERNIE (Zhang et al., 2019), K-BERT (Liu et al., 2020), DeBERTa (He et al., 2020), DABERT(Wang et al., 2022), DC-Match(Zou et al., 2022), DAFA(Song et al., 2022) and Large Language Models(LLMs) such as GPT3 (Brown et al., 2020), LLaMa (Touvron et al., 2023) have achieved superior results in many NLP tasks.

Although the pre-trained model shows a strong representation ability in sentence encoding, there are still some improvements for multi-head attention, which is used to improve the coding ability of the pre-trained model and improve the performance effect on downstream tasks, Such as 1) Syntax Bert(Bai et al., 2021) improves the model's understanding of text sentence structure by adding syntactic information to attention, 2) UER Bert(Xia et al., 2021) enhances the attention mechanism in the pre-trained model by introducing synonym information, and 3) SemBert(Zhang et al., 2020) improves the effect of text representation by integrating semantic role tagging and multi-label semantics into attention. In this work, we propose a scheme, which is used to provide commix different dimensional attention modes to capture the relationship between different components in a sentence, interactively augment and aggregate the four different attention modes (Commix Dimensional Attention) is used to improve the encoding ability of multi-head attention for text, and significantly improve the matching of short text, which can be easily combined with PLMs to stack additional improvements for text semantic matching.

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

237

238

239

241

242

243

244

245

246

247

248

249

250

251

253

254

255

256

257

259

260

261

262

263

264

3 Approach

We show the Commix Dimensional Attention Language Model in Figure 2. We take Bert as our base model, According to the findings in (Xia et al., 2021), the attention in the first layer transformer of Bert is broad and uninformed, so in order to improve the performance of the BERT as much as possible without adding extra parameters, We decided to apply our commix dimensional attention mechanism only in the first transformer layer of bert. Regarding this point, we also compared the results of changes in attention in each layer and changes in attention in all layers in the experiments. It also indirectly verified that adding only changes to attention in the first layer is optimal, please refer to the 4.2 for more details. Which consists of four parts under the augmentation-aggregation framework. Specifically, For every word embedding from **q** and **k**, we can obtain four matching scores using four different dimensional attention functions. Next, we augment the matching information along with words in **q**. We match two vectors inside each attention function interactively and then combine the matching information from all functions. The Multi-Layer Perception(MLP) is applied to fuse the matching information both in the self-interactive augmentation and normalized aggregation. Finally, that is, we obtain an aggregated commix dimensional attention result after different dimensional attention, self-interactive augmenta-



Figure 2: The overall architecture of the CDA-LM is shown in (a). The detailed structures of the Commix Dimensional Attention module, Self-interactive Augmentation module, and Normalized Aggregation module are shown in (b), (c), and (d), respectively.

tion, and normalized aggregation, which is used to replace the self-attention part in the first layer transformer of Bert.

3.1 Commix Dimensional Attention

266

267

269

270

272

273

274

275

281

284

287

291

292

296

297

In the commix dimensional attention module, we use four different dimensional attention functions to model the semantic relationship between sentence pairs from different perspectives. Note that this is in stark contrast to the single Attention used by default in the transformer. The four attentions are dot attention, additive attention, minus attention, and bilinear attention, respectively. The input of the commix dimensional attention module is a triple of $Q, K, V \in \mathbb{R}^{d_{seq}} \times d_v$, where d_v is the latent dimension, d_{seq} is the length of the utterance. We use q_i , k_i and v_i to denote the *i*-th dimension of Q, K, and V respectively. Four independent attention mechanisms compute potential relationships between Q, K, and V to measure their semantic interaction alignment.

3.1.1 Dot Attention

Dot attention is part of the commix dimensional attention module, which can directly compute correlations using matrix operations, and the computed scores are correlation weights. It is also the most commonly used attention mechanism in semantic correlation modeling. And it follows the standard dot-product attention that the transformer operates by default. The input of the dot attention module consists of queries and keys of dimension d_k , and values of dimension d_v . We compute the dot products of the query with all keys, and apply a softmax function to obtain the weights on the values. For the sake of simplicity, the formulations of BERT not be repeated here, please refer to (Devlin et al., 2018) for more details. We denote the output vector as:

$$\mathbf{s}_j^t = \mathbf{q}_j \odot \mathbf{k}_t$$
 (1a)

300

301

302

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

$$a_i^t = \frac{exp(s_i^t)}{\sum_{j=1}^N exp(s_j^t)}$$
(1b)

$$\mathbf{q}_t^d = \sum_{i=1}^N a_i^t \mathbf{v}_i \tag{1c}$$

where $\mathbf{q}_t^d \in R^{1 \times d_v}$ is the output of the *t*-th position obtained after the dot attention calculation and \odot is element-wise dot product.

3.1.2 Additive Attention

The second part of commix dimensional attention is the additive attention module, which is more inclined to capture aligned representations between sentence pairs from a global perspective since it applies a concatenated alignment of two vectors. Specifically, its input is the concatenation of two vectors, the output after the activation function represents the correlation between the vectors, and finally, the softmax function is applied to get the weights of the values. We denote the output vector as:

$$s_j^t = \tanh(\mathbf{W}_a([\mathbf{q}_j; \mathbf{k}_t])$$
 (2a) 32

$$a_i^t = \frac{exp(s_i^t)}{\sum_{j=1}^N exp(s_j^t)}$$
(2b) 325

$$\mathbf{q}_t^a = \sum_{i=1}^N a_i^t \mathbf{v}_i \tag{2c} 323$$

where $\mathbf{q}_t^a \in R^{1 \times d_v}$ is the output of the *t*-th position obtained after the additive attention calculation, and $\mathbf{W}_a \in R^{1 \times 2d_v}$ are weights of our model.

3.1.3 Minus Attention

327

329

330

331

333

334

335

339

340

341

342

343

344

345

354

360

The third part of commix dimensional attention is the minus attention module which captures and aggregates the different information between sentence pairs. The difference attention module adopts a subtraction-based cross-attention mechanism, which allows the model to pay attention to dissimilar parts between sentence pairs by elementwise subtraction as:

336
$$s_j^t = anh(\mathbf{W}_m(\mathbf{q}_j - \mathbf{k}_t))$$

$$a_i^t = \frac{exp(s_i^t)}{\sum_{j=1}^N exp(s_j^t)}$$
(3b)

$$\mathbf{q}_t^m = \sum_{i=1}^N a_i^t \mathbf{v}_i \tag{3c}$$

where $\mathbf{q}_t^m \in R^{1 \times d_v}$ is the output of the *t*-th position obtained after the minus attention calculation, and $\mathbf{W}_m \in R^{1 \times d_v}$ are weights of our model.

3.1.4 Bilinear Attention

The last part of commix dimensional attention is a bilinear attention module, which can learn a bilinear attention distribution of two vectors to seamlessly utilize the given sentence pair information. It models the bilinear interaction between two sets of input channels, facilitating the extraction of a joint representation of each pair of channels. Its calculation formula is as follows:

$$s_j^t = \mathbf{q}_j^T \mathbf{W}_b \mathbf{k}_t \tag{4a}$$

$$a_i^t = \frac{exp(s_i^t)}{\sum_{j=1}^N exp(s_j^t)}$$
(4b)

$$\mathbf{q}_t^b = \sum_{i=1}^N a_i^t \mathbf{v}_i \tag{4c}$$

where $\mathbf{q}_t^b \in R^{1 \times d_v}$ is the output of the *t*-th position obtained after the bilinear attention calculation and $\mathbf{W}_b \in R^{d_v \times d_v}$ are weights of our model.

3.2 Self-interactive Augmentation

Self-interactive augmentation is to fuse each word in the query vector in each attention function. For each position t, we splice the word representation h_t^k of v_t with its representation q_t^c of the correspond-361 ing attention, where v_t denotes the *t*-th dimension 362 of K, c = (a, b, d, m). which will better augment 363 each word representation with single attention to 364 capture the input information. then using a gating 365 structure to select the importance after splicing, and 366 then applying an MLP to fuse the representation 367 of each position more fully after the above opera-368 tions. After that, the output after self-interactive 369 augmentation is obtained. As shown below, This is 370 an example of our self-interactive augmentation of 371 additive attention: 372

$$\mathbf{x}_t^a = \begin{bmatrix} q_t^a, h_t^k \end{bmatrix}$$
(5a)

$$g_i = sigmoid\left(\mathbf{W}_g \mathbf{x}_t^a\right) \tag{5b}$$

$$\mathbf{x}_t^{a*} = g_i \odot \mathbf{x}_t^a \tag{5c}$$

$$\mathbf{h}_t^a = \tanh(\mathbf{W}_d \mathbf{x}_t^{a*} + b_a) \tag{5d}$$

For bilinear, dot, and minus attention, we will also get h_t^b , h_t^d , and h_t^m , respectively. Where $\mathbf{W}_g \in R^{1 \times 2d_v}$, $\mathbf{W}_d \in R^{d_v \times 2d_v}$, b_a are weights and bias of our model.

3.3 Normalized Aggregation

(3a)

Normalized aggregation is to fuse all the attention functions. We use a parameter z^c as an input to adaptively fuse four different attention mechanisms.

$$s_j = \tanh(\mathbf{W}_1 h_t^j + \mathbf{W}_2 \mathbf{z}^c) (j = a, b, d, m)$$

378 379

380

381

383

384

388

390

392

393

394

395

397

398

399

$$a_i = \frac{exp(s_i)}{\sum_{j=(a,b,d,m)} exp(s_j)}$$
(6b)

$$\mathbf{x}_t = \sum_{i=(a,b,d,m)} a_i \mathbf{h}_t^i \tag{6c}$$

Then, we input this X_t into an MLP neural network to fuse the information in different attention functions, and we will obtain different h_t^o for different positions in P.

$$\mathbf{h}_t^o = \tanh(\mathbf{W}_t \mathbf{x}_t + b_t) \tag{7}$$

Where $\mathbf{W}_t \in R^{d_v \times d_v}$, b_t are weights and bias of our model, respectively.

After aggregating the commix dimensional attention matching information, we will obtain the fused representation of the vectors at different positions for t from 1 to N.

$$A^{o} = (h_{1}^{o}, h_{2}^{o}, \dots, h_{N}^{o})$$
(8) 400

Table 1: The performance comparison of CDA-LM with other methods. We report Accuracy \times 100 on 6 GLUE datasets. Methods with \ddagger indicate the results from their papers, while methods with \ddagger indicate our implementation.

Method	Pre-trained	MRPC	QQP	MNLI-m/mm	QNLI	RTE	STS-B	Avg
BiMPM [†] (Wang et al., 2017)	×	79.6	85.0	72.3/72.1	81.4	56.4	-	-
CAFE [†] (Tay et al., 2017)	×	82.4	88.0	78.7/77.9	81.5	56.8	-	-
ESIM [†] (Chen et al., 2016)	×	80.3	88.2	-	80.5	-	-	-
Transformer [†] (Vaswani et al., 2017)	×	81.7	84.4	72.3/71.4	80.3	58.0	73.6	74.53
BiLSTM+ELMo+Attn [†] (Peters et al., 2018)	\checkmark	84.6	86.7	76.4/76.1	79.8	56.8	73.3	76.24
OpenAI GPT [†] (Radford et al., 2018)	\checkmark	82.3	70.2	82.1/81.4	87.4	56.0	80.0	77.06
UERBERT‡(Xia et al., 2021)	\checkmark	88.3	90.5	84.2/83.5	90.6	67.1	85.1	84.19
SemBERT [†] (Zhang et al., 2020)	\checkmark	88.2	90.3	84.4/84.0	90.9	69.3	87.3	84.90
BERT-base‡(Devlin et al., 2018)	\checkmark	87.2	89.0	84.3/83.7	90.4	66.4	85.8	83.83
RoBERTa‡(Liu et al., 2019b)	\checkmark	87.9	89.2	84.7/84.1	90.7	67.2	86.7	84.43
SyntaxBERT-base†(Bai et al., 2021)	\checkmark	89.2	89.6	84.9/84.6	91.1	68.9	88.1	85.20
CDA-LM-base [‡]	\checkmark	89.1	92.0	84.9/85.3	92.1	69.8	88.9	86.13
BERT-large‡(Devlin et al., 2018)	\checkmark	89.3	89.3	86.8/85.9	92.7	70.1	86.5	85.80
RoBERTa-large‡(Devlin et al., 2018)	\checkmark	90.4	89.4	86.8/86.1	92.7	72.3	87.5	86.32
SyntaxBERT-large†(Bai et al., 2021)	\checkmark	92.0	89.5	86.7/86.6	92.8	74.7	88.5	87.26
CDA-LM-large‡	\checkmark	91.9	92.3	87.3/87.4	95.2	75.7	89.8	88.59

 A^o is the final fused semantic feature and it will be propagated to the next computation flow.

Table 2: The performance comparison of CDA-LM with other methods on 4 popular datasets, including SNLI, Scitail(Sci), SICK, and TwitterURL(Twi).

Model	SNLI	Sci	SICK	Twi
ESIM [†] (Chen et al., 2016)	88.0	70.6	-	-
CAFE†(Tay et al., 2017)	88.5	83.3	72.3	-
CSRAN [†] (Tay et al., 2018)	88.7	86.7	-	84.0
BERT-base‡(Devlin et al., 2018)	90.7	91.8	87.2	84.8
RoBERTa-base‡(Liu et al., 2019b)	90.9	92.3	87.9	85.9
UERBERT‡(Xia et al., 2021)	90.8	92.2	87.8	86.2
SemBERT†(Zhang et al., 2020)	90.9	92.5	87.9	86.8
MT-DNN-base†(Liu et al., 2019a)	91.1	94.1	-	-
SyntaxBERT-base†(Bai et al., 2021)	91.0	92.7	88.5	87.3
CDA-LM-base‡	91.8	94.0	89.2	88.2
BERT-large‡(Devlin et al., 2018)	91.0	94.4	91.1	91.5
RoBERTa-large‡(Liu et al., 2019b)	91.2	94.5	91.2	91.9
SyntaxBERT-large†(Bai et al., 2021)	91.3	94.7	91.4	92.1
CDA-LM-large‡	92.1	95.5	92.9	92.8

4 Experiment

The datasets, baselines, and all details of our experiments are shown in Appendix A.3.

4.1 Results

In the experiments, we replace the original attention module with our CDA mechanism in the BERT Table 3: The performance comparison of CDA-LM-large with LLaMa family and GPT3 on several datasets. including SNLI, Scitail(Sci), SICK, and TwitterURL(Twi).

Model	SNLI	Sci	SICK	Twi
LLaMA-7B‡(Touvron et al., 2023)	73.1	75.2	68.4	64.7
LLaMa-13B‡(Touvron et al., 2023)	78.5	83.2	80.6	83.1
GPT3‡(Brown et al., 2020)	84.3	90.1	88.7	85.8
CDA-LM-large‡	92.1	95.5	92.9	92.8

model.

Firstly, we fine-tune our model on 6 GLUE datasets. Table 1 shows the performance of CDA-LM compared with some other baseline models. It can be seen that the effect of non-pre-trained models is significantly worse than pre-trained model has more data from the learning corpus and a pow-erful information extraction ability. The performance of our CDA-based BERT-base and BERT-large model improves the original BERT models by 2.2% and 2.7%, respectively. Moreover, our model also outperforms SyntaxBert (which is the state-of-the-art external knowledge-based model) by 0.9% on BERT-base and 1.3% on Bert-large, respectively.

Secondly, to verify the overall performance of our method, we also conduct experiments on four other popular datasets. The results are shown in Table 2, CDA-LM outperforms vanilla Bert and 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

501

452

453

454



Figure 3: Stability experiments on QQP (a), SNLI (a), QNLI (b), Scitail(b) datasets.



Figure 4: layer-by-layer experiments on MRPC and QQP datasets

other competitive models on almost all datasets. In addition, the amount of data in Scitail is relatively small, which makes the variance of the model prediction results larger. However, CDA-LM still shows very competitive performance on Scitail, which also shows that our method can make up for the lack of generalization ability with fewer parameters by endowing Bert with subtle difference awareness.

Overall, consistent conclusions can be drawn from these results. Compared with previous work, our method shows very competitive performance in judging semantic similarity, and the experimental results also confirm our idea.

4.1.1 Vs LLMs

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

We also compared the semantic understanding abil-444 ity of the Large Language Models(LLMs) with 445 our model. Considering that the answers provided 446 by the Large Language Models(LLMs) involve 447 content analysis and generation, in order to make 448 the comparison results more comparative, we per-449 formed the following operations on the experiment. 450 For the convenience of statistics, we have extracted 451

partial data randomly from several datasets as test data. In addition, We have made the following settings for the prompt: Please provide the similarity between the following two sentences. If similar, provide 1; if dissimilar, provide 0.

From the table 3, we can see that although LLMs have strong comprehension and generation abilities, there are still some shortcomings in dealing with semantic matching tasks, mainly due to the inability to understand prompts well. For example, for some sentence pairs, providing some unwanted answers cannot directly provide 1 or 0, especially in large-scale data processing, which can generate many answers that cannot be programmed and batch processed in downstream tasks. In addition, there may be some deviations in the understanding of sentences, leading to incorrect judgments. Moreover, some sentence pairs generate very unstable answers. For example, after the model's judgment, sometimes it gives a result of 1, but when executed again later, the model gives a result of 0. These also indicate the instability of large language models in semantic matching tasks.

4.1.2 Stability Analysis

We also performed extensive experiments on QQP, SNLI, QNLI, and Scitail datasets to explore the stability of our method. To minimize the impact of randomness in Bert's training, performance levels were averaged over 10 different runs on the development set. The performance distribution box diagram is shown in Figure 3. The median and average levels of our model exceeded the ordinary Bert on all four datasets, and the performance fluctuation range of our method was within $\pm 1\%$ of the average level, which indicates that our method has better stability than Bert on different data distributions.

4.2 layer-by-layer Analysis

Regarding why we only applied commix attention to the first layer transformer in BERT, in addition to being inspired by the article in UERBERT (Xia et al., 2021), we also made changes in each layer and also in all layers, as shown in the figure 4. We applied commix attention to the attention of transformers in each layer of BERT, and at the end, we also made changes to transformers in all layers of BERT, selecting the MRPC and QQP datasets, The experimental results are basically consistent with the conclusion in (Xia et al., 2021), that is, the information understanding capability of the first layer



Figure 5: Distribution of Dot attention (a), Additive attention (b), Minus attention (c) and Bilinear attention (d).

transformer in BERT is the worst, and modifying the commix attention in this layer is the most effective, better than applying commix attention in all layers, and also better than applying commix attention in all layers.

4.3 Attention Distribution

502

503

505

507

To visually demonstrate the impact of different at-508 tention functions inside commix dimensional atten-509 tion on the interactive alignment of sentence pairs, 510 we show the weight distribution of four kinds of 511 attention in the figure 5. We can observe that the 512 word-pair information in the sentence pairs con-513 cerned with different attention functions is incon-514 sistent. First, in Figure (a), Dot attention can pay at-515 tention to the same words and semantically related 516 words in sentence pairs, but it is heavily influenced 517 by the same words in sentence pairs. It focuses too 518 much on the shallow features of the same text and 519 ignores the deep semantic association of the dif-520 ferent words between "software" and "hardware". 521 This shows that using Dot attention alone may lead to wrong predictions. Secondly, in Figure (b), we can see that the distribution of attention weights is 524 more uniform, because the calculation method of 525 additive attention tends to fuse the two signals, so 526 it pays attention to different word pairs of software 528 and hardware to a certain extent, but the interaction weights still small. Next, in Figure (c), it can be observed that Minus attention explicitly pays 530 attention to the difference between "software" and "hardware", and its attention weight is the largest 532 among all word pairs. This is because subtractive 533 attention uses element-wise subtraction to compare 534 the differences between sentence pairs. The greater the difference between word pairs, the greater their 536 weight. Therefore, it can also be complementary 537 to Dot attention. Finally, in Figure (d), the atten-538 tion weights in bilinear attention focus on the same words, which indicates that bilinear attention tends

to focus on the same parts of sentence pairs, and this mechanism is beneficial for capturing sentence pairs' commonality. In summary, different attention focus on different word pairs in sentence pairs. Intuitively, commix dimensional attention can effectively combine the alignment relationships of multiple views in sentence pairs to generate vectors that better describe the matching details of sentence pairs. 541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

561

563

564

565

567

568

569

570

571

572

573

574

575

576

577

578

579

The more result of our additional experiments is shown in Appendix A.

5 Conclusion

In this paper, we propose a novel Commix Dimensional Attention (CDA) mechanism to improve the large-scale pre-trained model, such as BERT, and RoBERTa for the semantic matching task. The commix dimensional attention module fully exploits complementary and complex relationships between sentences compared to the single attentionbased model. Moreover, the self-interactive augmentation enables better interaction between each attention function and its input, enhancing the representation ability of each attention mechanism. Furthermore, the proposed adaptive aggregation module with normalized aggregation mechanisms can effectively fuse the key features and filter out the unrelated features produced by the commix dimensional attention module for semantic understanding. Extensive experiments on 6 GLUE benchmark datasets, as well as 4 other commonly used semantic understanding datasets, verify that the proposed CDA-LM achieves remarkable performance improvements over the original single attention mechanism-based BERT model as well as other state-of-the-art semantic understanding models. Since the CDA mechanism is a universal transformation mechanism for transformers, it is expected to be applied to other large-scale pre-trained models in the future.

Limitations

580

This work has the following limitations: (1) The proposed method is based on the introduction of 582 multiple attention functions. Since the introduced attention functions have not been pre-trained, if they are not fine-tuned on the labeled dataset, errors 585 586 may be introduced and propagated to the decision model, resulting in label prediction errors. (2) We 588 initially demonstrated that external structures can be combined with BERT to improve performance on various SSM tasks. We are also interested in 590 trying to combine it with other PLMs. However, due to computational resource constraints, we did 592 not conduct more experiments on other PLMs. (3) 594 Introducing an extrinsic structure significantly improves the generalization ability of PLMs in fewshot scenarios, but a deeper understanding of why this is the case is still lacking. This may inspire 597 better methods to exploit pre-trained models.

References

599

606

607

610

611

613

614

615

616

617

618

619

623

625

- Jiangang Bai, Yujing Wang, Yiren Chen, Yaming Yang, Jing Bai, Jing Yu, and Yunhai Tong. 2021. Syntaxbert: Improving pre-trained transformers with syntax trees. *arXiv preprint arXiv:2103.04350*.
- Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Qian Chen, Xiaodan Zhu, Zhenhua Ling, Si Wei, Hui Jiang, and Diana Inkpen. 2016. Enhanced lstm for natural language inference. *arXiv preprint arXiv:1609.06038*.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. 2019. What does bert look at? an analysis of bert's attention. *arXiv preprint arXiv:1906.04341*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Xu Han, Zhengyan Zhang, Ning Ding, Yuxian Gu, Xiao Liu, Yuqi Huo, Jiezhong Qiu, Yuan Yao, Ao Zhang, Liang Zhang, et al. 2021. Pre-trained models: Past, present and future. *AI Open*, 2:225–250.

Yaru Hao, Li Dong, Furu Wei, and Ke Xu. 2020. Self-attention attribution: Interpreting information interactions inside transformer. *arXiv preprint arXiv:2004.11207*, 2. 630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770– 778.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735– 1780.
- Tushar Khot, Ashish Sabharwal, and Peter Clark. 2018. Scitail: A textual entailment dataset from science question answering. In *Thirty-Second AAAI Confer*ence on Artificial Intelligence.
- Wuwei Lan, Siyu Qiu, Hua He, and Wei Xu. 2017. A continuously growing dataset of sentential paraphrases. *arXiv preprint arXiv:1708.00391*.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*.
- Yongjie Lin, Yi Chern Tan, and Robert Frank. 2019. Open sesame: getting inside bert's linguistic knowledge. *arXiv preprint arXiv:1906.01698*.
- Weijie Liu, Peng Zhou, Zhe Zhao, Zhiruo Wang, Qi Ju, Haotang Deng, and Ping Wang. 2020. K-bert: Enabling language representation with knowledge graph. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 2901–2908.
- Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. 2019a. Multi-task deep neural networks for natural language understanding. *arXiv preprint arXiv:1901.11504*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019b. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, Roberto Zamparelli, et al. 2014. A sick cure for the evaluation of compositional distributional semantic models. In *Lrec*, pages 216–223. Reykjavik.

- 710 711 712 714 716 721 722 723 727 729 731

- Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. arXiv preprint arXiv:1802.05365.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding with unsupervised learning.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. arXiv preprint arXiv:1910.01108.
- Jian Song, Di Liang, Rumei Li, Yuntao Li, Sirui Wang, Minlong Peng, Wei Wu, and Yongxin Yu. 2022. Improving semantic matching through dependencyenhanced pre-trained model with adaptive fusion. arXiv preprint arXiv:2210.08471.
- Chuanqi Tan, Furu Wei, Wenhui Wang, Weifeng Lv, and Ming Zhou. 2018. Multiway attention networks for modeling sentence pairs. In IJCAI, pages 4411-4417.
- Yi Tay, Anh Tuan Luu, Aston Zhang, Shuohang Wang, and Siu Cheung Hui. 2019. Compositional deattention networks. In Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc.
- Yi Tay, Luu Anh Tuan, and Siu Cheung Hui. 2017. A compare-propagate architecture with alignment factorization for natural language inference. arXiv preprint arXiv:1801.00102, 78:154.
- Yi Tay, Luu Anh Tuan, and Siu Cheung Hui. 2018. Costack residual affinity networks with multi-level attention refinement for matching text sequences. arXiv preprint arXiv:1810.02938.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Advances in neural information processing systems, pages 5998-6008.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. arXiv preprint arXiv:1804.07461.
- Sirui Wang, Di Liang, Jian Song, Yuntao Li, and Wei Wu. 2022. Dabert: Dual attention enhanced bert for semantic matching. arXiv preprint arXiv:2210.03454.

Zhiguo Wang, Wael Hamza, and Radu Florian. 2017. Bilateral multi-perspective matching for natural language sentences. arXiv preprint arXiv:1702.03814.

736

738

739

740

741

742

743

744

745

746

747

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

766

- Tingyu Xia, Yue Wang, Yuan Tian, and Yi Chang. 2021. Using prior knowledge to guide bert's attention in semantic textual matching tasks. In Proceedings of the Web Conference 2021, pages 2466–2475.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. Advances in neural information processing systems, 32.
- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model. arXiv preprint arXiv:2210.02414.
- Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. Ernie: Enhanced language representation with informative entities. arXiv preprint arXiv:1905.07129.
- Zhuosheng Zhang, Yuwei Wu, Hai Zhao, Zuchao Li, Shuailiang Zhang, Xi Zhou, and Xiang Zhou. 2020. Semantics-aware bert for language understanding. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, pages 9628–9635.
- Yicheng Zou, Hongwei Liu, Tao Gui, Junzhe Wang, Qi Zhang, Meng Tang, Haixiang Li, and Daniell Wang. 2022. Divide and conquer: Text semantic matching with disentangled keywords and intents. In Findings of the Association for Computational Linguistics: ACL 2022, pages 3622-3632, Dublin, Ireland. Association for Computational Linguistics.

A Appendix

769

770

771

772

774

775

777

778

779

780

790

791

794

796

801

803

804

806

808

810

811

812

814

815

816

817

818

A.1 The Result of Additional Experiments

A.1.1 Ablation Study

To evaluate the contribution of each component in our method, we conduct ablation experiments on the QQP datasets based on BERT. The experimental results are shown in table 4.

The commix dimensional attention module consists of four core components that use different attention functions to model the correlation between sentence pairs. We want to know if each component is useful for the sentence-matching task. First, after removing dot attention, the performance of the model drops by 2.3%. While dot attention can capture the dynamic alignment relationship between word pairs, which is crucial for semantic matching tasks. Then, When additive attention and bilinear attention are removed respectively, the performance of the model on the datasets drops to 90.4% and 90.2%, respectively. They are significantly smaller than Dot attention, which indicates that these two kinds of attention are weaker than dot attention in distinguishing sentence pair relations. Finally, after removing minus attention from the model, the performance dropped by 2.1%. The different information can further describe the interaction between words and can provide more fine-grained comparison information for the pretrained model so that the model can obtain a better representation. The above experiments show that the performance drops when the sub-module is removed, which demonstrates the effectiveness of the internal components of the commix dimensional attention module.

Next, in the aggregation module, we also conduct multiple experiments to verify the effect of augmentation and aggregation of multiple matching features. On the QQP datasets, we first remove the self-interactive augmentation module, and the performance drops to 90.8%. Since selfinteractive augmentation can capture interactions between multiple signals, this interaction information is crucial for fusing multi-source vectors. Second, after removing the normalized aggregation module, we only integrate multiple signals by simple averaging. The accuracy drops to 90.5%, which proves that dynamic aggregation according to different weights can further improve the performance of the model. Finally, when we remove the both augmentation module and aggregation module and use simple averaging instead, the performance drops sharply to 89.6%, which is the largest drop among all ablation components. This suggests that while commix dimensional attention mechanisms are crucial for judging sentence pair relations, hard-integrating multi-attention mechanisms without interactive augmentation into PLMs may destroy their pre-existing knowledge, while selfinteractive augmentation and soft aggregation can better enhance and aggregate multiple attention. 819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

Table 4: Results of component ablation experiment.

Model	Dev	Test
CDA-LM	92.6	92.0
w/o Dot attention	90.4	89.7
w/o Additive attention	91.1	90.4
w/o Minus attention	90.7	89.9
w/o Bilinear attention	91.3	90.2
w/o Self-interactive augmentation	91.9	90.8
w/o Normalized aggregation	91.6	90.5
w/o Augmentation and Aggregation	90.3	89.6

Finally, to explore whether any two kinds of attention can serve as complementary cues, we aggregate the four kinds of attention in pairs, and the performance improvement over baseline BERT is shown in Figure 6. First, we can find that after the aggregation of different kinds of attention, the performance of the baseline based on a single attention mechanism is improved. Second, The fusion of dot attention and minus attention achieves the best complementary performance improvement among all aggregations, and the combination of dot attention and other attentions significantly outperforms other attention pairs, which reflects that dot attention contributes more to the text matching task than other attentions. Finally, it is worth noting that other types of attention aggregation minus attention can also achieve better results. This may be because minus attention can capture the different information in sentence pairs, and can intuitively reflect the differences between sentence pairs.

Overall, due to the efficient combination of each component, CDA-LM can adaptively fuse vectors generated from different attentions from multiple perspectives into a pre-trained model and leverage its powerful contextual representation to better infer semantics.

Table 5: The robustness experiments results of CDA-LM and other models. The data transformation methods we utilized mainly include SwapAnt(SA), NumWord(NW), AddSent(AS), InsertAdv(IA), AppendLrr(AL), AddPunc(AP), BackTrans(BT), TwitterType(TT), SwapNameEnt(SN), SwapSyn-WordNet(SW)

Model		Q	uora			SNLI				
	SA	NW	IA	AL	BT	AS	SA	TT	SN	SW
ESIM [†] (Chen et al., 2016)	-	-	-	-	-	64.00	84.22	78.32	53.76	65.38
DistilBERT [†] (Sanh et al., 2019)	42.24	56.85	83.10	84.09	83.20	-	-	-	-	-
BERT [†] (Devlin et al., 2018)	48.58	56.96	86.32	85.48	83.42	79.66	94.84	83.56	50.45	76.42
ALBERT†(Lan et al., 2019)	51.08	55.24	81.87	78.94	82.37	45.17	96.37	81.62	57.66	74.93
SyntaxBERT†(Bai et al., 2021)	49.30	56.37	86.43	84.62	84.19	78.63	95.31	86.91	58.26	76.90
CDA-LM‡	55.93	63.26	87.75	85.08	87.99	81.56	97.35	85.64	60.62	81.23
Model	MNLI-m/mm									
Wodel	AS	S	А	А	P	Т	Т	S	N	SW
BERT‡(Devlin et al., 2018)	55.32/55.25	52.76/55.69		82.30/82.31		77.08	/77.22	51.97	/51.84	76.41/77.05
ALBERT†(Lan et al., 2019)	53.09/53.58	50.25	/50.20	83.98	/83.68	77.98	/78.03	56.43	/50.03	76.63/77.43
SyntaxBERT†(Bai et al., 2021)	54.92/54.63	53.54	/54.73	80.01	/79.71	75.46	/74.93	57.11	/51.95	78.57/79.31
CDA-LM‡	60.75/59.82	58.33	/60.83	83.63	/83.59	78.24	/78.36	60.77	/60.35	82.58/83.21



Figure 6: The influence of different attention integration methods on the QQP test set. The baseline model is BERT-base. Att_{dot} indicates Dot Attention, Att_{add} indicates Additive Attention, Att_{min} indicates Minus Attention and Att_{bil} indicates Bilinear Attention.

A.1.2 Robustness Performance Test

855

857

858

862

870

871

To examine the performance of CDA-LM and competition models in terms of their ability to capture subtle differences in sentence pairs. We performed robustness tests on three widely studied datasets. Table 5 lists the accuracy of the 6 models on the three data sets. We can observe that SwapAnt leads to maximum performance degradation, which indicates that the model cannot handle the semantic contradiction expressed by antonyms (non-explicit negations) between sentence pairs. The model performance on NumWord drops to 63.26% in Quora datasets, because it requires the model to capture subtle numerical differences for correct language reasoning. Meanwhile, ESIM performed worst. The results reflect that the pre-trained mechanism benefits from abundant external resources and provides better generalization ability than the denovo

training model. The performance of the improved pre-trained model SyntaxBERT is better than that of the original Bert model, which reflects that sufficient pre-trained corpus and appropriate external knowledge fusion strategy are helpful to improve the generalization performance of the model. On TwitterType and AddPunc, the performance of CDA-LM is lower than that of AlBERT but still better than that of Bert, which may be related to the pre-trained corpus and training mechanism. In the other 8 conversions, CDA-LM can show attention to subtle differences and obtain better performance. 873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

A.2 Model Parameter Analysis

Our model is based on the fusion of four different kinds of attention, meanwhile, attention modules are often used to explore the interpretability of the model (Clark et al., 2019; Hao et al., 2020; Lin et al., 2019), in order to prove the impact of the increase of parameter quantity on the model effect, and the impact of multiple fusion of a single attention and the effect of multiple fusion models of different attention, as shown in Table 6, we have conducted experiments on four separate attention on the QQP datasets, and for these four separate attention, 4 times the number of attention are used to participate in the calculation to achieve the same level of parameter quantity as our CDA-LM. The experimental results show that:

 The effect of a single kind of attention model is basically the same as that of a model with a 903single kind of attention expanded by four times,904which means that the increase of the number of905parameters of the model caused by a simple four-906fold increase in the number of single attention907does not have a significant improvement on the908performance of the model.

909

910

911

912

913

914

915

Our CDA-LM compared to the model with a single kind of attention expanded by four times have a significant improvement of the result, This shows that the effective fusion of different attachments does have complementary advantages in semantic extraction.

Model	Dev	Test
CDA-LM	92.6	92.0
Dot attention $\times 1$	90.4	89.1
Dot attention $\times 4$	91.2	89.5
Additive attention $\times 1$	90.7	88.9
Additive attention $\times 4$	91.8	89.3
Minus attention $\times 1$	91.9	89.9
Minus attention \times 4	90.1	90.2
Bilinear attention $\times 1$	89.6	88.8
Bilinear attention $\times 4$	90.3	89.2

Table 6: Results of model parameter experiment.

A.3 Case Study

In order to intuitively understand how CDA-LM 916 works, we use the three cases in Table 7 for quali-917 tative analysis. First, although S1 and S2 are literally similar in the first example, they express two 919 completely different semantics due to the subtle difference the phrases bring to "eat fruit" and "eat 921 early". The non-pre-trained language model ESIM 922 is difficult to capture the semantic conflict caused by the different words in case 1, so ESIM gives wrong prediction results. The pre-trained language model BERT can identify semantic differences in case 1 and give correct predictions with the help of strong contextual representation capabilities. It is worth noting that the similarity of Bert's predicted 929 sentence pairs is 46.32%, while that of CDA-LM is only 1.87%. Second, in case 2, "from 70 to 60" and 931 "from 60 to 50" in sentence pairs express different 932 semantics, but they are mainly caused by numerical 933 differences. Although BERT identified the correct label in case 1 by a small margin, in case 2 it was 935 unable to capture numerically induced differences 936 and gave wrong predictions because it requires the 937 model to capture subtle numerical differences for correct language reasoning. Finally, our model

made correct predictions in all of the above cases. Since CDA-LM models sentence pairs from multiple perspectives, it can pay attention to the small differences in sentence pairs, and adaptively aggregate multi-source information in the aggregation module to better identify the semantics within sentence pairs' differences. At the same time, we can observe that ESIM performs the worst, and the results reflect that the pre-trained mechanism benefits from abundant external resources and provides better generalization ability than the denovo training model. And our BERT-based improved model CDA-LM outperforms the original BERT model, reflecting that a reasonable structure improvement and an effective aggregation strategy can further improve the model's generalization performance.

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

A.4 Implementation Details of Our Experiments

Implementation Details CDA-LM is based on BERT-base and BERT-large. We set the number of both self-attention layers and heads as 12, and the dimension of embedding vectors as 768. The total number of trainable parameters of both the original BERTbase and our proposed model is the same 110M For distinct targets, besides, our hyper-parameters are different. We use AdamW(Loshchilov and Hutter, 2017) in the BERT and set the learning rate in $\{1e^{-5}, 2e^{-5}, 3e^{-5},$ $8e^{-6}$. As for the learning rate decay, we use a warmup(He et al., 2016) of 0.1 and L2 weight decay of 0.01. Furthermore, we set the epoch to 5 and the batch size is selected in $\{16, 32, 64\}$. We also set dropout at 0.1-0.3. To prevent gradient explosion, we set gradient clipping in $\{7.5, 10.0, 15.0\}$. All the experiments are conducted by Tesla V100 and PyTorch platform. In addition, to ensure that the experimental results are statistically significant, we conduct each experiment five times and report the average results.

A.5 Datasets Statistics

The statistics of all 10 datasets are shown in Table 8.

A.5.1 GLUE Datasets

We experimented with 6 datasets of the GLUE¹ datasets (Wang et al., 2018): MRPC, QQP, STS-B, MNLI-m/mm, QNLI, and RTE, The following is a detailed introduction to these datasets.

¹https://huggingface.co/datasets/glue

Case	ESIM	BERT	CDA-LM
S1: Can eat only fruit for dinner lead to weight loss?S2: Does ate dinner earlier in the evening help with weight loss?	similarity:92.76%	similarity:46.32%	similarity:1.87%
	label:1	label:0	label:0
S1: How do girls lose weight from 70 to 60 ?S2: How should I lose weight from 60 to 50 ?	similarity:99.51%	similarity:72.66%	similarity:12.06%
	label:1	label:1	label:0
S1: What skills do I need to learn to be a successful hardware engineer?S2: What should I do to become a successful software engineer?	similarity:99.99%	similarity:99.26%	similarity:18.63%
	label:1	label:1	label:0

Table 7: The example sentence pairs of our cases. Red and Blue are different phrases in sentence pair.

- **MRPC** is a dataset that automatically extracts sentence pairs from online news sources and manually annotates whether the sentences in sentence pairs are semantically equivalent. The task is to determine whether there are two categories of interpretation: interpretation or not interpretation.
- **QQP** comes from the famous community Q&A website quora. Its goal is to predict which of the provided question pairs contains two questions with the same meaning.
- **STS-B** is a collection of sentence pairs extracted from news headlines, video titles, image titles, and natural language inference data. Each pair is annotated by humans, and its similarity score is 0-5. The task is to predict these similarity scores, which is essentially a regression problem, but it can still be classified into five text classification tasks of sentence pairs.
- **MNLI-m/mm** is a crowd-sourced collection of sentence pairs annotated with textual entailment information. Given the promise statement and hypothesis statement, the task is to predict whether the premise statement contains assumptions (entailment), conflicts with assumptions (contradiction), or neither (neutral).

QNLI is a question and answer data set composed of a question paragraph pair, in which the paragraph is from Wikipedia, and a sentence in the paragraph contains the answer to the question. The task is to judge whether the question and sentence (sentence, a sentence in a Wikipedia paragraph) contain, contain and do not contain, and classify them.

• **RTE** is a series of datasets from the annual text implication challenge. These data samples are constructed from news and Wikipedia. All these data are converted into two categories. For the data of three categories, neutral and contradiction are converted into not implication in order to

Datasets	#Train	#Dev	#Test	#Class
MRPC	3669	409	1380	2
QQP	363871	1501	390965	2
MNLI-m/mm	392703	9816/9833	9797/9848	3
QNLI	104744	40432	5464	2
RTE	2491	5462	3001	2
STS-B	5749	1500	1379	2
SNLI	549367	9842	9824	3
SICK	4439	495	4906	3
Scitail	23596	1304	2126	2
TwitterURL	42200	3000	9324	2

Table 8: The statistics of all 10 datasets.

maintain consistency.

A.5.2 Other Datasets

We also experimented with 4 other popular datasets: SNLI², Scitail³, SICK⁴ and TwitterURL⁵. The following is an introduction to these 4 datasets.

- **SNLI**(Bowman et al., 2015) is a dataset used for classification (or natural language inference). The task is to determine whether two sequences entail, contradict or are mutually neutral.
- **Scitail**(Khot et al., 2018) is an entailment dataset created from multiple-choice science exams and web sentences. Each question and the correct answer choice are converted into an assertive statement to form the hypothesis.
- **SICK**(Marelli et al., 2014) is a dataset for semantic textual similarity estimation. The task is to assign a similarity score to each sentence pair.
- **TwitterURL**(Lan et al., 2017) is a collection of sentence-level paraphrases from Twitter by linking tweets through shared URLs. Its goal is to discriminate between duplicates and not.

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1018

1019

1020

1021

1022

1023

1024

1025

1026

987

989

²https://nlp.stanford.edu/projects/snli/

³https://allenai.org/data/scitail

⁴http://marcobaroni.org/composes/sick.html

⁵https://github.com/lanwuwei/Twitter-URL-Corpus

1048

1050

1051

1052

1053

1054

1055

1057

1058

1059

1060

1061

1062

1063

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1078

1079

1080

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1094

1095

1096

1097

1098

A.6 Baselines

To evaluate the effectiveness of our proposed CDA-LM in SSM, we mainly introduce BERT (Devlin et al., 2018), SemBERT (Zhang et al., 2020), SyntaxBERT(Liu et al., 2020), UERBERT (Xia et al., 2021) and multiple PLMs (Radford et al., 2018; Devlin et al., 2018) for comparison. Moreover, we also selected several competitive no pre-trained models as baselines, such as ESIM (Chen et al., 2016), Transformer (Vaswani et al., 2017), etc (Hochreiter and Schmidhuber, 1997; Wang et al., 2017; Tay et al., 2017). Besides, We also compared with Large Language Models, such as GPT3, and the LLaMa family.

- **BIMPM** is proposed in (Wang et al., 2017) proposes a Bilateral Multi-Perspective Matching (Bilateral Multi-Perspective Matching, BiMPM) model for sentence matching.
- **CAFE** (Tay et al., 2017) introduces a new architecture where alignment pairs are compared, compressed, and then propagated to upper layers for enhanced representation learning. And then it adopts factorization layers for efficient and expressive compression of alignment vectors into scalar features, which are then used to augment the base word representations.
- ESIM (Chen et al., 2016) is a model that combines BiLSTM and attention proving that the sequential inference model based on chained LSSM can outperform previous complex structures. It further achieved new SOTA performances.
- **CSRAN** (Tay et al., 2018) is a deep architecture, involving stacked recurrent encoders. CSRAN incorporates two novel components to take advantage of the stacked architecture. It first introduces a new bidirectional alignment mechanism that learns affinity weights by fusing sequence pairs across stacked hierarchies. And then it leverages a multi-level attention refinement component between stacked recurrent layers.
 - **Transformer** (Vaswani et al., 2017) uses the attention mechanism to reduce the distance between any two positions in the sequence to a constant. It is not a sequential structure similar to RNN, so it has better parallelism.
- ELMO (Peters et al., 2018) adopts a typical twostage process. The first stage is pre-training using a language model; The second stage is to extract the word embedding of each layer of the network corresponding to the word from the pre-training

network and add it to the downstream task as a new feature. It can solve the problem of polysemy of the previous language model because the generated word vector is changed according to the change of the specific use context. 1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

- **GPT** (Radford et al., 2018) is a semi-supervised learning method that uses a large amount of unlabeled data to let the model learn "common sense" to alleviate the problem of insufficient labeled information. The specific method is to pre-train the model Pretrain with unlabeled data before training Fine-tune for labeled data, and ensure that the two kinds of training have the same network structure.
- **BERT** (Devlin et al., 2018) Given that our model implements based on BERT, we naturally compare it with vanilla BERT without prior knowledge. We adopt the configuration of Google's BERT-base in our experiments.
- UERBERT (Xia et al., 2021) conducted lots of experiments to analyze which kind of external knowledge BERT has already known, and directly injected the synonym knowledge into BERT without fine-tuning.
- SemBERT(Zhang et al., 2020) incorporates contextual semantics from pre-trained semantic role labeling and is capable of explicitly absorbing contextual semantics over a BERT backbone. SemBERT keeps the convenient usability of its BERT precursor in a light fine-tuning way without substantial task-specific modifications.
- **Syntax-BERT** (Liu et al., 2020) is a framework that integrates the syntax trees into transformerbased models. Unlike us, it explicitly injected syntactic knowledge into checkpoints of models.
- **MT-DNN**(Liu et al., 2019a) not only leverages large amounts of cross-task data but also benefits from a regularization effect that leads to more general representations to help adapt to new tasks and domains. MT-DNN extends the model by incorporating a pre-trained bidirectional transformer language model.
- GPT3 (Brown et al., 2020) As the successor of GPT-2, GPT-3 has become the largest language model, further expanding the parameter space (175 billion vs. 1.5 billion) and data size (45 TB vs. 40 GB). The model requires no fine-tuning to formulate downstream tasks and has excellent performance in zero-shot and few-shot settings. Based on the multi-task generalization capabilities of GPT-2, GPT-3 has achieved good results

- 1150on many new tasks, including mathematical ad-
dition, news article generation, vocabulary inter-
pretation, and code writing. As the number of
parameters increases, the model becomes more
powerful.
- 1155 • LLaMa(Touvron et al., 2023) trained on a new mix of publicly available data. LLaMa is a collec-1156 tion of pretrained and fine-tuned large language 1157 models (LLMs) ranging in scale from 7 billion to 1158 70 billion parameters. Testing conducted to date 1159 1160 has been in English and has not and could not cover all scenarios. Therefore, before deploying 1161 any applications of LLaMa, developers should 1162 perform safety testing and tuning tailored to their 1163 specific applications of the model to facilitate the 1164 safe deployment of LLaMa. 1165