DiFiNet: Boundary-Aware Semantic Differentiation and Filtration Network for Nested Named Entity Recognition

Anonymous ACL submission

Abstract

001 Nested Named Entity Recognition (Nested NER) entails identifying and classifying en-003 tity spans within the text, including the detection of named entities that are embedded within external entities. Prior approaches primarily employ span-based techniques, utilizing the power of exhaustive searches to ad-007 dress the challenge of overlapping entities. Nonetheless, these methods often grapple with the absence of explicit guidance for boundary detection, resulting insensitivity in discerning minor variations within nested spans. To this end, we propose a Boundary-aware Semantic Differentiation and Filtration Network (DiFiNet) tailored for nested NER. Specifically, DiFiNet leverages a biaffine attention mech-017 anism to generate a span representation matrix. This matrix undergoes further refinement through a self-adaptive semantic differentiation module, specifically engineered to discern semantic variances across spans. Furthermore, DiFiNet integrates a boundary filtration module, designed to mitigate the impact of nonentity noise by leveraging semantic relations among spans. Extensive experiments on three benchmark datasets demonstrate our model yields a new state-of-the-art performance¹.

1 Introduction

Named Entity Recognition (NER) involves the utilization of computer-assisted techniques to identify and extract entities and corresponding semantic types (Lample et al., 2016a), including *person* (PER), *location* (LOC), *geo-political entity* (GPE), and others. NER plays a crucial role in facilitating various downstream tasks such as relation extraction (Tang et al., 2022), event extraction (Yang and Mitchell, 2016; Sha et al., 2018; Yao et al., 2021), and sentiment analysis (Zhao et al., 2021).

conventional approaches have primarily focused on identifying non-nested entities (Chiu and

Another tornado hits Geneva, near	the Alabama - Florida line, said	
	GPE GPE	
	LOC	
	GPE	

Figure 1: A sample sentence from ACE Corpus containing nested entities.

Nichols, 2016; Lample et al., 2016b; Ma and Hovy, 2016), a trend largely attributed to the constraints of corpus annotations that emphasize flat entity structures. However, the complex nature of natural language frequently features nested named entities, with studies revealing that approximately 30% of sentences in ACE04 and ACE05 datasets contain such structures (Finkel and Manning, 2009; Kati-yar and Cardie, 2018). The prevalence of nested structures underscores the need for efficient models adept at handling such linguistic complexities.

In response to this challenge, recent years have witnessed a burgeoning interest in nested NER (Ju et al., 2018; Luo et al., 2024; Wang et al., 2020; Tan et al., 2021a). Among the emerging strategies, span-based models stand out as prominent approaches and have set new benchmarks in the field (Tan et al., 2020; Wang and Lu, 2020; Zhong and Chen, 2021; Zhu and Li, 2022). These models excel by leveraging exhaustive search techniques to systematically identify all possible spans, thereby capturing the full spectrum of nested structures.

Despite the success of span-based methods, they often struggle to fully utilize the rich semantics within spans due to the absence of explicit guidance for boundary detection. Previous research indicates that span-based models usually encounter confusion when dealing with nested entities characterized by a high degree of token overlap (Tan et al., 2020; Wan et al., 2022). To illustrate, consider the sentence taken from the ACE05 dataset in Figure 1, entities like "*the Alabama-Florida line*", "*Florida*", and "*Alabama*", as well as non-entity spans such as "*Alabama-Florida line*" or "*the Alabama-Florida*"

073

074

041

¹The source code is anonymized online at: https:// anonymous.4open.science/r/DiFiNet-00CE/

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

165

166

167

168

169

170

171

172

173

share a significant number of tokens, blurring the semantic distinction between entity and non-entity spans. Besides, those low-quality candidate spans, particularly long entities, incur significant computational costs (Tan et al., 2020; Shen et al., 2021; Tan et al., 2021b) due to the extensive array of potential spans evaluated during training, inevitably limited by practical constraints.

075

076

077

084

087

090

096

100

101

102

103

104

105

106

107

109

110

111

112

113

114

115

116

To address this issue, we propose a Boundaryaware Semantic Differentiation and Filtration network (DiFiNet) explicitly incorporating semantic differences between nested spans as input features. By leveraging gradient back-propagation, DiFiNet learns appropriate internal representations to augment the distinction among nested entities within the span semantic representation space, enhancing its ability to discern boundaries and accurately classify nested named entities. Specifically, DiFiNet integrates BERT and a biaffine attention mechanism to construct a matrix of span semantic representations, followed by a self-adaptive semantic differentiation module to transform span representations into semantic differences across spans. Additionally, to alleviate the influence of low-quality candidate spans within the matrix, DiFiNet integrates a boundary filtration module. This module serves to model the interaction among spans, effectively reducing noise, with a specific focus on distinguishing semantically similar entity and non-entity spans. Our main contributions are summarized as follows:

> • We tackle the challenge of nested named entity recognition from a novel perspective by explicitly enhancing boundary supervision to address the issue of boundary insensitivity within nested entities.

• Building upon our perspective, We propose a novel end-to-end framework which effectively captures subtle semantic variations between entity and non-entity spans. This framework is engineered to precisely detect entity boundaries via both self-adaptive semantic differentiation and boundary filtration module.

Extensive experiments on the ACE04, ACE05, and GENIA datasets indicate that DiFiNet outperforms existing state-of-the-art models in the nested NER task. Further ablation studies validate the contribution of each module within our framework.

2 Related Work

Nested Named Entity Recognition is a task in Natural Language Processing (NLP) that involves identifying and classifying named entities within text data, where entities can have complex and overlapping structures. One approach to tackle this task is the hypergraph method, originally proposed by Lu and Roth (2015). This method maps the nested entity structures to sub-graphs in a hyper-graph and performs classification on them. Several extensions have been developed based on this method (Muis and Lu, 2017; Katiyar and Cardie, 2018).

Another approach is the hierarchical method introduced by Ju et al. (2018), which divides entities into different levels, where each deeper level represents a higher level of entity specificity. Following this paradigm, Wang et al. (2020) designed a pyramid sequence labeling framework using convolutional neural networks to extract entities from bottom to top. Shibuya and Hovy (2020) explored suboptimal path decoding to progressively extract entities hierarchically, and Wang et al. (2021) further improved it by excluding the influence of the optimal path. However, both hyper-graph and hierarchical methods suffer from high complexity when dealing with complex nested entities.

In contrast, Seq2Seq methods offer a simpler end-to-end approach, typically utilizing LSTM-CRF (Straková et al., 2019) or BART (Yan et al., 2021) to predict the label of each position. Zhang et al. (2022) improved Seq2Seq methods by adopting intra-entity and inter-entity de-confounding data augmentation techniques. Shen et al. (2023b) designed a dual-slot multi-prompt template with a position slot for locating and a type slot for typing, respectively. Nevertheless, when faced with highly complex nesting structures, these methods may encounter long-distance dependency problems, resulting in cascading errors.

To address the aforementioned challenges in nested NER, Sohrab and Miwa (2018) proposed a span-based method that treats the nested NER task as span prediction problems. This approach involves predicting potential entity spans for each token, followed by filtering and merging these spans to obtain the final nested entities. Building upon Sohrab's work, several works have made advancements to the span-based method by incorporating graph structure (Wan et al., 2022), valuable span patterns (Shen et al., 2021; Tan et al., 2021a) and attention mechanism (Yu et al., 2020; Xu et al.,

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

267

268

223

224

225

226

227

2021; Zheng et al., 2023) to achieve state-of-the-art performance. For example, Yu et al. (2020) proposed a biaffine attention mechanism to enhance the interaction between the start and end tokens and assigned scores to each span. When constructing the span-based contrastive loss function, Zhang et al. (2023) utilizes concatenation to generate span representations. Shen et al. (2023a) redefined NER by modeling it as a boundary-denoising diffusion process. This approach generates named entities by refining and clarifying noisy spans.

174

175

176

177

178

179

180

181

183

184

187

188

189

191

192

193

194

195

196

197

199

204

206

209

210

211

212

213

214

215

216

217

218

219

However, span-based models typically utilize pooling (Eberts and Ulges, 2020; Shen et al., 2021; Li et al., 2021), concatenation (Li et al., 2021; Tan et al., 2020; Zheng et al., 2023) or integration (Zhu and Li, 2022; Yuan et al., 2022; Shen et al., 2023a) techniques to generate span representations from token representations. However, this approach often leads to generating semantically similar representations for highly overlapping spans. As a result, effectively capturing the subtle semantic nuances within individual spans becomes challenging.

To mitigate the boundary insensitivity issue, we propose to explicitly incorporate span semantic difference features into nested NER task. This allows the model to learn more robust span representations by capturing the nuanced semantic variations between entity and non-entity spans.

3 Our Approach

In this section, we introduce the details of our framework as shown in Figure 2. We first formulate the task definition of nested NER as follow,

Nested NER as boundary detection In the context of nested NER, the task involves analyzing an input sentence denoted as X = $\{x_1, x_2, \ldots, x_n\}$ to identify and classify potential entities according to a predefined set of entity types $T = \{t_1, t_2, \ldots, t_k\}$. Typically, an entity can be represented by a triplet (s_i, e_i, t_i) , where s_i and e_i denote the starting and ending position of the entity, respectively, and $t_i \in T$ represents the assigned entity type. This structured representation allows for the precise localization of entity boundaries within the sentence. In a sentence with *n* tokens, there are a total of n(n + 1)/2 valid spans.

3.1 Span Semantic Encoder

Given a sentence $X = \{x_1, x_2, ..., x_n\}$, we first utilize a pre-trained BERT model (Devlin et al., 2019) to vectorize each token x_i , resulting in tokenlevel feature representations denoted as $\mathbf{H}_{enc} = \{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n \mid \mathbf{h}_i \in \mathcal{R}^{d \times 1}\}$, where *d* is the embedding dimension, and *n* denotes the number of tokens within the sentence.

We then design two feed-forward neural networks (FNNs) to map the tokens and obtain the semantic representation vectors for the start and end tokens \mathbf{h}_s , $\mathbf{h}_e \in \mathcal{R}^{l \times h}$ of a span, where l represents the sentence length and h denotes the hidden dimension. Subsequently, a biaffine model is employed to combine the start and end token representation, and the width representation of the span $(\mathbf{w}_{ij} \in \mathcal{R}^c)$ to construct the span representation matrix $\mathbf{M}^0 \in \mathcal{R}^{l \times l \times f}$, where f corresponds to the number of biaffine features.

For each span S_{ij} , spanning from the *i*-th token to the *j*-th token, its vector \mathbf{M}_{ij}^0 is computed as:

$$\mathbf{h}_{s} = \operatorname{GELU} \left(\mathbf{H}_{enc} \mathbf{W}_{s} \right); \mathbf{h}_{e} = \operatorname{GELU} \left(\mathbf{H}_{enc} \mathbf{W}_{e} \right), \\ \mathbf{M}_{ij}^{0} = \left(\mathbf{h}_{s}[i] \oplus \mathbf{h}_{e}[j] \oplus \mathbf{w}_{ij} \right) \mathbf{W} + \mathbf{h}_{s}[i] \mathbf{U} \mathbf{h}_{e}[j]^{T},$$
 (1)

where $\mathbf{W}_s, \mathbf{W}_e \in \mathcal{R}^{h \times h}, \mathbf{W} \in \mathcal{R}^{(2h+c) \times r}$, and $\mathbf{U} \in \mathcal{R}^{h \times r \times h}$ are learnable parameters. The feature size of biaffine model is denoted by $r. \oplus$ denotes concatenation, and GELU refers to the *gelu* activation function. It is worth noting that when \mathbf{M}_{ij}^0 is situated off the diagonal of the \mathbf{M}^0 matrix, the span representation \mathbf{S}_{ij} exhibits two distinct forms, symmetrically arranged along the diagonal.

3.2 Self-adaptive Semantic Differentiation Module

To effectively capture the semantic differences between spans, we propose the **Self-adaptive Differentiation** operator (SAD), inspired by computer vision techniques such as the Roberts cross operator (Roberts and Lawrence, 1965). The SAD operator addresses the rigid nature of traditional gradient operators by adapting its differentiation template to the local semantic context of each span, bolstering the capability of handling subtle variations between semantically similar entity and non-entity spans.

The SAD operator functions in two primary phases: *the masking phase* and *the differentiation phase*. During *the masking phase*, a learnable convolutional kernel assesses local semantic regions, generating a mask matrix $mask_{x_0}$ that highlights the most pertinent neighboring spans for semantic differentiation, formulated as:

$$I_{x_0} = \underset{i \in R \setminus \{x_0\}}{\operatorname{arg max}} (\operatorname{LN}(\operatorname{Conv}(\mathbf{M}^0_{x_0}))), \quad (2)$$



Figure 2: An overview of DiFiNet with two-layer structure SDM. \bigcirc denoted the concatenate operation. \oplus denoted the element-wise addition operation. \otimes denoted Hadamard product operation. LN denotes LayerNorm layer.

269

270

274 275

277

278 279

280

282

283

284

287

289

291

$$SADBlock(*) = GELU(LN(SAD(*))),$$

$$\mathbf{M}_{lr}^{1} = SADBlock(SADBlock(\mathbf{M}^{0})),$$
(5)

 $\mathbf{mask}_{x_0}[i] = \begin{cases} 1 & if \ i = I_{x_0} or \ i = x_0 \\ 0 & others \end{cases}, \quad (3)$

where $Conv(\mathbf{M}_{x_0}^0)$ represents the convolution op-

eration on R centered around x_0 , changing the

channel number from f to the number of spans

in R. LN denotes the layer normalization opera-

tion, and arg max represents the position of the

The differentiation phase employs the mask ma-

trix to apply self-adaptive weights to the span rep-

resentations. This is achieved by element-wise

multiplication of $mask_{x_0}$ with a fixed weight ma-

trix \mathbf{w}_{f} , enabling nuanced semantic differentiation

 $SAD(x_0) = \sum_{x_n \in B} \mathbf{mask}_{x_0} \cdot \mathbf{w}_f \cdot \mathbf{M}_{x_n}^0,$

where $\mathbf{M}_{x_n}^0$ denotes the span semantic matrix at

Integrated within the Self-adaptive Semantic

Differentiation Module (SDM), the SAD operator

underpins two SAD Blocks in each layer, designed

to capture both first-order and second-order seman-

tic differences between spans, denoted as:

position x_n and \mathbf{w}_f is the fixed weight matrix.

maximum score in R except for x_0 .

tailored to each span's context:

where $lr \in \{0, 1, ..., N\}$ and N + 1 is the number of layers in the model. For the sake of simplicity, the equations presented do not include the residual connections in SAD Blocks.

293

294

295

296

297

298

300

301

302

303

304

305

306

307

308

310

311

312

To enable back-propagation of gradients in the SAD operator, which contains a non-differentiable Argmax operation, we employ the Gumbel softmax estimator (Jang et al., 2016). Additionally, to ensure consistency in the differentiated objects, the fixed weights of the SAD operators in SDM have opposite signs. The weight matrix \mathbf{w}_f is used for the first SAD operator, while the second SAD operator adopts the matrix \mathbf{w}'_f whose elements are the negations of \mathbf{w}_f :

Subsequently, the SDM processes semantic difference features, aligning them within a standardized semantic framework. These processed features are then integrated using a linear layer, which consolidates the individual semantic distinctions into a comprehensive span boundary matrix \mathbf{M}_{fuse}^2 :

$$\mathbf{M}_{lr}^{2} = \operatorname{Conv}_{1 \times 1}(\mathbf{M}_{lr}^{1}),$$
$$\mathbf{M}_{fuse}^{2} = \mathbf{W}_{fuse}(\underbrace{\mathbf{M}_{0}^{2} \oplus \dots \oplus \mathbf{M}_{lr}^{2}}_{N+1}) + \mathbf{B}_{fuse}, \quad (7)$$

361

362

363

364

365

366

367

369

370

371

372

373

374

375

376

377

378

380

381

383

384

385

386

387

389

390

391

392

393

394

395

where \mathbf{W}_{fuse} is the weight matrix and \mathbf{B}_{fuse} is the bias term of the linear layer. The convolution operator $\operatorname{Conv}_{1\times 1}$ denotes a 2D convolution operation with 1×1 kernel, and \oplus indicates concatenation.

3.3 Boundary Filtration Module

318

319

320

322

324

327

332

333

335

336

338

339

341

342

343

345

347

351

355

357

The introduction of low-quality information, particularly in distinguishing non-entity spans, presents challenges in the SDM, potentially leading to an increase in false positives. To mitigate this, we introduce the **Boundary Filtration Module** (BFM), designed to reduce the impact of irrelevant span by extracting and utilizing interactions between spans, which aids in clarifying entity boundaries.

The BFM utilizes a structured methodology incorporating a top-down pathway for semantic interaction extraction and a bottom-up approach for detail restoration, complemented by lateral connections for comprehensive span relationship analysis. The top-down pathway employs a series of convolution blocks that apply Layer Normalization and the GELU activation function to refine span features systematically:

$$ConvBlock(*) = GELU(LN(Conv(*))),$$
$$\mathbf{M}_{g}^{0} = ConvBlock(\mathbf{M}^{0}), \qquad (8)$$
$$\mathbf{M}_{g}^{i} = ConvBlock(\mathbf{M}_{g}^{i-1}),$$

where $i \in \{1, 2, ..., n\}$ and n + 1 represents the number of convolution blocks.

The bottom-up pathway, in contrast, aims to restore finer details from higher-layer features through up-sampling, using nearest neighbor techniques to retain critical relational information. This is synchronized with lateral connections to merge features from different layers effectively, thereby avoiding loss of detail and preventing checkerboard artifacts typically associated with interpolation methods:

$$\mathbf{M}_{g}^{i-1'} = \text{upSample}(\mathbf{M}_{g}^{i}) + \mathbf{M}_{g}^{i-1},$$

$$\mathbf{M}_{g} = \text{Conv}(\text{upSample}(\mathbf{M}_{g}^{0'})).$$
 (9)

3.4 Span Semantic Decoder

In order to preserve the complete semantic information of span, we incorporate \mathbf{M}^0 as residual to \mathbf{M}_{fuse}^2 . The composite matrix then undergoes linear decoding to yield prediction logits:

$$\mathbf{p} = \sigma(\mathbf{W}_p(\mathbf{M}^0 \oplus \mathbf{M}_{fuse}^2 \oplus \mathbf{M}_g) + \mathbf{B}_p), \quad (10)$$

where $\mathbf{p} \in \mathcal{R}^{l \times l \times t}$, $\mathbf{W}_p \in \mathcal{R}^{d \times t}$, $\mathbf{B}_p \in \mathcal{R}^t$. \mathbf{W}_p and \mathbf{B}_p are trainable parameters. σ denotes *Sigmod* activation function.

3.5 Training and Inference

Training We minimize the following binary cross-entropy loss function:

$$\mathcal{L} = -\sum_{0 \le i,j < l} \mathbf{y}_{ij} \log(\mathbf{p}_{ij}) + (1 - \mathbf{y}_{ij}) \log(1 - \mathbf{p}_{ij}),$$
(11)

where \mathbf{y}_{ij} is the ground truth entity type. To accommodate DiFiNet's architecture, which does not distinguish between the matrix halves during training, we incorporate errors from both the upper and lower triangular sections of \mathbf{p}_{ij} and \mathbf{p}_{ji} , aligning with the symmetric nature of entity representation.

Inference For entity prediction, we average the values from the upper and lower sections of **p** to ensure consistent decoding:

$$\mathbf{p}_{ij}' = (\mathbf{p}_{ij} + \mathbf{p}_{ji})/2. \tag{12}$$

Following Yu et al. (2020), we first eliminate spans deemed non-entities (those with all probabilities below 0.5), then rank the remaining spans by their highest probability. Spans are selected sequentially; any span conflicting with previously chosen spans in terms of boundaries is omitted, maintaining clear entity demarcation.

4 Experiment

4.1 Datasets

We evaluate our model on three commonly used nested NER datasets: ACE04², ACE05³, and GE-NIA ⁴. For the ACE datasets, we use the data preprocessing code released by Yan et al. (2023) and split the data into training, validation, and test sets by 8:1:1. For the GENIA dataset, we follow Li et al. (2022) to categorize entities into five types and split data into train, dev and test sets by 8:1:1. See Appendix A for detailed information of datasets.

4.2 Baselines

To evaluate the performance of the proposed model, we compare it with the following models on three datasets: Biaffine (Yu et al., 2020), Second-Best (Wang et al., 2021), Seq2Seq (Yan et al., 2021), Sequence2Set (Tan et al., 2021a), De-bias(Zhang

²https://catalog.ldc.upenn.edu/ LDC2005T09

³https://catalog.ldc.upenn.edu/ LDC2006T0

⁴http://www.geniaproject.org/ genia-corpus

Table 1: The performance of various models on the ACE04, ACE05, and GENIA datasets is presented in Table 1. The "Encoder" column indicates the pre-trained models utilized by each model for the ACE datasets, while all models employed BioBERT-Base for the GENIA dataset. \dagger signifies that the models were reproduced using the same pre-processed data and publicly available code. The best results are highlighted in **bold** font. The subscript denotes the standard deviation, providing a measure of result variability (e.g., 88.48_{23} indicates a value of 88.48 ± 0.23).

Models	Encodor	ACE04			ACE05			GENIA		
widdels	Liicodei	Р	R	F1	Р	R	F1	Р	R	F1
Biaffine (2020)	BERT-base	87.30	86.00	86.70	85.20	85.60	85.40	81.8	79.30	80.50
Second-Best (2021)	BERT-base	86.42	85.71	86.06	83.95	84.67	84.30	79.20	78.16	78.63
Locate-and-Label (2021)	BERT-base	87.44	87.38	87.41	86.09	87.27	86.67	80.19	80.89	80.54
Seq2Seq (2021)	BART-large (2020)	87.27	86.41	86.84	83.16	86.38	84.74	78.57	79.30	78.93
Sequence2Set (2021a)	BERT-large	88.46	86.10	87.26	87.48	86.63	87.05	82.30	78.70	80.40
Span-Graph (2022)	BERT-base	86.70	85.93	86.31	84.37	85.87	85.11	77.92	80.74	79.30
De-bias (2022)	-bias (2022) T5-base (2020)		84.54	85.44	83.31	86.56	84.90	81.04	77.21	79.08
BS (2022) † RoBERTa-base		87.32	86.84	87.08	86.58	87.84	87.20	82.53	78.69	80.56
Triaffine (2022) †	BERT-large	87.13	87.68	87.40	86.70	86.94	86.82	80.42	82.06	81.23
W2NER (2022) †	BERT-large	87.19	87.72	87.45	85.77	87.80	86.76	83.10	79.76	81.39
ICR (2023)	BERT-large	-	-	-	87.11	87.14	87.12	79.02	80.68	79.87
BINDER(2023) †	BERT-large	87.34	88.30	87.81	87.41	88.34	87.87	81.69	80.85	81.26
DiffusionNER(2023a) †	sionNER(2023a) † BERT-large		87.52	87.42	85.04	88.42	86.70	81.85	79.59	80.70
PromptNER(2023b) † BERT-large		87.02	88.03	87.52	86.01	88.12	87.05	-	-	-
CNNNER (2023) †	RoBERTa-base	87.33	87.29	87.31	86.70	88.16	87.42	83.19	79.70	81.40
DiFiNet	RoBERTa-base BERT-large	88.57 88.64	88.43 88.32	88.45 ₁₄ 88.48 ₂₃	89.16 88.62	88.74 88.17	88.94 ₃₈ 88.39 ₃₁	83.01	80.80	81.87 ₁₉

et al., 2022), W2NER (Li et al., 2022), Locate-and-Label (Shen et al., 2021), BS (Zhu and Li, 2022), Triaffine(Yuan et al., 2022), Span-Graph (Wan et al., 2022), ICR(Zheng et al., 2023), BINDER (Zhang et al., 2023), CNNNER (Yan et al., 2023), DiffusionNER(Shen et al., 2023a) and Prompt-NER(Shen et al., 2023b). See Appendix B and C for further elaboration on baseline models and implementation details of DiFiNet, respectively.

4.3 Main Results

398

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

our evaluation employs three key metrics: **Precision**, **Recall**, and **F1-score**, to assess the performance of the models. We adopt strict evaluation criteria, whereby precise matches in both entity boundaries and categories are required for correct recognition. To validate the consistency and reliability of our findings, we conducted five separate trials, each initialized with distinct random seeds, and then proceeded to statistical analysis on the collected F1 scores. Specifically, we applied the T-test at a 5% significance level to determine the statistical significance of the differences observed between experimental outcomes.

Table 1 presents a comprehensive performance of DiFiNet and baseline models on ACE04, ACE05, and GENIA datasets for NER. Across all three NER datasets, DiFiNet consistently outperforms the baseline models. Notably, with RoBERTa-base as the underlying pre-trained model, DiFiNet secures an increase of +1.14% in F1-score on ACE04 and +1.52% in F1 on ACE05 compared to existing models. Similarly, when leveraging BERT-large as the pre-trained backbone, DiFiNet attains enhancements of +0.67% F1 on ACE04 and +0.52% F1 on ACE05. Additionally, DiFiNet exhibits an improvement of +0.47% F1 on the GENIA dataset. It is essential to highlight that the marginal gains on the GENIA dataset might stem from its significantly lower frequency of nested entities (18.41%) compared to ACE04 (45.68%) and ACE05 (39.11%), as shown in Table 6. These results underscore the superior performance of DiFiNet in addressing the complexities of nested NER.

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

4.4 Ablation Studies

Table 2 reports the F1 score results of the DiFiNet and its variants. The variations explored include disabling the Self-adaptive Semantic Differentiation Module (SDM) (w/o SDM), removing the Boundary Filtration Module (BFM) (w/o BFM), and excluding both (w/o SDM, BFM). Additionally, to gauge the impact of the self-adaptive mechanism, we examine a configuration without the selfadaptive mask (w/o Self-adaptive mask). Each variant demonstrates a drop in F1 score compared to DiFiNet model, highlighting the individual and collective importance of these modules in enhanc-

483

484

n = 0	87.48	-0.97	87.43
n = 2	87.64	-0.81	88.24
SAD Block \times 1	87.59	-0.86	88.18
SAD Block \times 3	87.26	-1.19	87.84
ing nested NER performance of the number of SAD Plant	ormanc ty of the s the nu	e. Fur model mber o	thermo to var f SDM
the number of SAD BI	OCKS W1	thin ea	ch SDI
and the number of co	nvoluti	on bloc	cks wit
BFM. The adjustments	s are ma	de whi	le maiı
other settings at their	optima	l levels	to iso

Settings

w/o Self-adaptive mask

DiFiNet

w/o SDM

w/o BFM

N = 0

N = 2

w/o SDM, BFM

Table 2: Ablation experiment results (RoBERTa-base as the pre-trained language model). Δ denotes the performance drops (F1 Score) under different experimental conditions compared to our proposed model.

F1

87.43

87.78

87.09

87.53

87.60

87.63

ACE04

88.45

Δ

-1.02

-0.67

-1.36

-0.92

-0.85

-0.82

ACE05

88.94

Δ

-1.15

-0.73

-1.81

-0.62

-1.83

-0.93

-1.51

-0.70

-0.76

-1.10

F1

87.79

88.21

87.13

88.32

87.11

88.01

ing nested NER performance. Furthermore, we
scrutinize the sensitivity of the model to various hy-
perparameters, such as the number of SDM layers,
the number of SAD Blocks within each SDM layer,
and the number of convolution blocks within the
BFM. The adjustments are made while maintaining
other settings at their optimal levels to isolate the
effects of each parameter.

Our findings indicate the following: (1) Necessity of SDM and BFM: The elimination of either the SDM, the BFM, or both significantly diminishes the model's effectiveness. Such a reduction underscores the essential roles that these modules play in identifying semantic variances across spans and in bolstering the model's ability to detect boundaries; (2) Adaptive Sampling Benefits: Adaptive sampling within the differentiation process improves performance, indicating limitations in static approaches for handling complex nested entity structures; (3) SDM Layer Impact: Additional SDM layers do not guarantee improved outcomes, suggesting an optimal level of model complexity that avoids unnecessary noise; (4) BFM Convolution Blocks: Excessive convolution blocks in BFM don't lead to better results and may remove essential information, indicating a balance is needed; (5) Optimization with SAD Blocks: The model performs best with two SAD blocks, showing that this balance effectively captures semantic differences without overcomplicating the model. Overall, the experiments validate the importance of each proposed module in optimizing model performance.

Table 3: Entity length-wise results on ACE05 dataset. Entities are divided into six groups based on their lengths. The % column represents the proportion of entities in each length range out of the total number, rounded to two decimal places.

		w/o S	DM and	BFM		DiF	liNet
Len.	%	Р	R	F1	Р	R	F1
[1, 4)	87.58	87.67	89.09	88.38	89.27	89.69	89.48 (+1.1)
[4, 7)	7.55	84.39	85.47	84.93	86.12	86.47	86.29 (+1.36)
[7, 10)	2.32	68.00	70.83	69.39	81.17	82.28	81.72 (+12.33)
[10,13)	1.00	72.22	83.87	77.61	72.97	87.10	79.41(+1.80)
[13,16)	0.55	56.53	76.47	65.00	82.25	77.47	79.79 (+14.79)
$[16,+\infty)$	1.00	56.10	74.19	63.89	68.94	73.63	71.74 (+8.85)

Table 4: Results on CoNLL03 dataset. All models utilize BERT-large as a pretrain encoder, and all results are from their respective original papers.

Models	CoNLL03				
Widdels	Р	R	F1		
W2NER	92.71	93.44	92.07		
DiffusionNER	92.99	92.56	92.78		
PromptNER	92.48	92.33	92.41		
BINDER	93.08	93.57	93.33		
DiFiNet	93.84	93.60	93.72		

4.5 **Performance on Long Entities**

Within NER tasks, the identification of long entities poses substantial challenges, notably due to a higher likelihood of encompassing nested structures, which exacerbates boundary insensitivity issues. Additionally, the accurate recognition of long entities represents a long-tail challenge (Wan et al., 2022), making their detection particularly complex. 485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

507

508

509

510

511

In Table 3, we present a comprehensive analysis of the impact of SDM and BFM on the recognition of entities with varying lengths. The results demonstrate that DiFiNet not only maintains but also enhances F1 scores for longer entities, with a remarkable improvement of +14.79% F1 for entities between 13 to 16 words in length and 8.85% improvement on entities over the length of 16. This significant enhancement is credited to DiFiNet's adeptness at discerning intricate semantic variances among overlapping and extended spans, facilitated by span semantic divergence features. By capitalizing on these features, DiFiNet substantially improves its handling and representation of long entities, which are typically characterized by more complex configurations and multiple nested spans.

4.6 Performance on Flat Entities

To evaluate the performance of our model on flat NER, we compared it against four leading state-of-

Table 5: Case Study on ACE05. The labels in the lower right corner indicate the entity type, while the superscripts indicate the nesting level. [The candidate entity]_T denotes predicted with incorrect type T. {^{m1}The candidate entity^{m1}} represents the missed ground true entities whose number m1.

	Sentence 1
Ground True / DiFiNet	$[^{1}$ These Iraqis ¹] _{PER} were rifling [¹ a home of [² a senior member of [³ the Mukhabarat ³] _{ORG} , [³ [⁴ Saddam ⁴] _{PER} 's dreaded secret police ³] _{ORG} ²] _{PER} ¹] _{FAC} .
CNNNER	$\begin{bmatrix} {}^{1}\text{These Iraqis}^{1} \end{bmatrix}_{\text{PER}} \text{ were rifling } \{ {}^{m1}\text{a home of } \{ {}^{m2}\text{a senior member of } [{}^{1}\text{the Mukhabarat}^{1}]_{\text{ORG}}, \\ [{}^{1}[{}^{2}\text{Saddam}^{2}]_{\text{PER}} \text{'s dreaded secret police}^{1}]_{\text{ORG}} {}^{m2} \} {}^{m1} \}.$
	Sentence 2
Ground True / DiFiNet	from $[^{1}$ the $[^{2}$ cnn ² $]_{ORG}$ center in $[^{2}$ atlanta ² $]_{GPE}$ $^{1}]_{FAC}$, $[^{1}i^{1}]_{PER}$ ' m $[^{1}$ fred fred $^{1}]_{PER}$.
CNNNER	from $[^{1}$ the $[^{2}$ cnn ² $]_{ORG}$ center in $[^{2}$ atlanta ² $]_{GPE}$ $^{1}]_{FAC}$, $[^{1}i^{1}]_{PER}$, m { m^{1} fred fred m^{1} }.
	Sentence 3
Ground True / DiFiNet	But [¹ neighboring Malaysia ¹] _{GPE} 's success in integrating [¹ [² Russian ²] _{GPE} MiG-29s ¹] _{VEH} and [¹ [² American ²] _{GPE} [2F/A-18 ²] _{VEH} Hornets ¹] _{VEH} persuaded [¹ them ¹] _{PER} otherwise, [Sudarsono] _{PER} said.
CNNNER	But [¹ neighboring Malaysia ¹] _{GPE} 's success in integrating [¹ [² Russian ²] _{GPE} MiG-29s ¹] _{VEH} and [¹ [² American ²] _{GPE} [2F/A-18 ²] _{VEH} Hornets ¹] _{VEH} persuaded [¹ them ¹] _{GPE} otherwise, [Sudarsono] _{PER} said.

the-art models on CoNLL03 dataset ⁵, W2NER (Li et al., 2022), DiffusionNER (Shen et al., 2023a), PromptNER (Shen et al., 2023b), and BINDER (Zhang et al., 2023). Table 4 shows that our model outperforms these benchmarks, particularly in precision metrics, achieving precision score of 93.84. This performance indicates that our model's explicit guidance for boundary detection not only aids nested entity recognition but also significantly enhances flat entity identification.

5 Case Study

512

513

514

515

516

518

519

520

522

523

524

525

526

527

529

532

533

536

537

538

540

Table 5 shows a case study conducted on ACE05 to compare DiFiNet with CNNNER (Yan et al., 2023). The first observation highlights that DiFiNet demonstrates superior ability to identify nested long entities due to its proficiency in detecting subtle distinctions between spans. The second sample demonstrates that DiFiNet excels in recognizing entities not encountered during training, leveraging semantic differences between spans. For instance, in the absence of training data for the boundary word "fred", it becomes challenging for the model to identify it based solely on span representation. However, by drawing guidance from the semantic difference between "i am fred fred" and "fred fred", DiFiNet can recognize the pattern of "*i am [name]*" in context, facilitating the accurate identification of the entity "fred fred". Furthermore, DiFiNet exhibits advantages in resolving ambiguous entity

references. By leveraging the semantic difference between "*persuaded them otherwise*" and "*them*", DiFiNet effectively recognizes the pattern of "*persuaded [person] otherwise*" and appropriately classifies "*them*" as PER. However, due to CNNNER lacking awareness of subtle semantic differences, it fails to correctly identify all entities in three examples. We provide extended case studies in Appendix D to further illustrate DiFiNet's ability to capture subtle semantic differences between spans. 541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

568

6 Conclusion

This paper proposes a Boundary-aware Semantic Differentiation and Filtration Network (DiFiNet) to effectively address the issue of boundary insensitivity in nested named entity recognition tasks. DiFiNet introduces the self-adaptive semantic differentiation module to capture semantic difference information between spans and incorporates the boundary filtration module to reduce noise from non-entity spans and enhance the differences of boundary semantics between spans. Experimental results demonstrate that DiFiNet achieves superior performance compared to existing approaches on three benchmark datasets. Ablation experiments and case studies further validate the effectiveness of the proposed model. Looking ahead, we aim to extend utilization of boundary information in tasks such as event extraction and relation extraction.

⁵https://www.clips.uantwerpen.be/ conll2003/ner/

661

662

663

664

665

669

569 Limitations

We discuss here the limitations of the method in 570 this paper. First, this method still needs to tra-571 verse all spans, bringing high computational costs. 572 Second, since the biaffine model encodes spans as continuous entities, it results in the prediction of only contiguous entities. Therefore, this method 575 has limited applicability for noncontinuous entity 576 recognition tasks. Finally, effectively integrating 577 multi-level span semantic difference information is a promising direction for optimization. 579

80 Ethics Statement

582

583

584

588

591

592

594

596

598

600

607

608

609

610

611

613

To ensure ethical considerations, we will provide a detailed description as follows:

- 1. All of the datasets used are collected and annotated in previous studies. The use of these datasets in our work does not involve any interaction or collection of individual privacy data.
- Our work focuses on methodology studies and experiments. The results and models in our paper will not be used to harm or deceive any individuals or groups.
- 3. There are no potential conflicts of interest or ethical issues regarding financial support in the sponsors and funds of our research work.

Acknowledgement

We would like to thank the anonymous reviewers for their helpful discussion and feedback. This work was supported by (ANONYMIZED FOR DOUBLE BLIND REVIEW).

References

- Jason PC Chiu and Eric Nichols. 2016. Named entity recognition with bidirectional lstm-cnns. *Transactions of the association for computational linguistics*, 4:357–370.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

- Markus Eberts and Adrian Ulges. 2020. Span-based joint entity and relation extraction with transformer pre-training. In *the 24th European Conference on Artificial Intelligence*, pages 2006–2013, Online. IOS Press.
- Jenny Rose Finkel and Christopher D Manning. 2009. Nested named entity recognition. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 141–150, Singapore. Association for Computational Linguistics.
- Eric Jang, Shixiang Gu, and Ben Poole. 2016. Categorical reparameterization with gumbel-softmax. In *International Conference on Learning Representations*, pages 1–13, San Juan, Puerto Rico. OpenReview.net.
- Meizhi Ju, Makoto Miwa, and Sophia Ananiadou. 2018. A neural layered model for nested named entity recognition. In *The 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, volume 1, pages 1446–1459, New Orleans, Louisiana. Association for Computational Linguistics.
- Arzoo Katiyar and Claire Cardie. 2018. Nested named entity recognition revisited. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 861–871, New Orleans, Louisiana. Association for Computational Linguistics.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016a. Neural architectures for named entity recognition. In 2016 Conference of the North American Chapter of the Association for Computational Linguistics, pages 260–270, San Diego, USA. Association for Computational Linguistics.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016b. Neural architectures for named entity recognition. In Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 260–270.
- Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. BioBERT: A pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics (Oxford, England)*, 36(4):1234–1240.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

- 670 671 672 679
- 694 702 705 709 710 711 712 713 714 715 716 717 718 719 720
- 721
- 725

- Fei Li, ZhiChao Lin, Meishan Zhang, and Donghong Ji. 2021. A span-based model for joint overlapped and discontinuous named entity recognition. In The 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, volume 1, pages 4814-4828, online. Association for Computational Linguistics.
- Jingye Li, Hao Fei, Jiang Liu, Shengqiong Wu, Meishan Zhang, Chong Teng, Donghong Ji, and Fei Li. 2022. Unified named entity recognition as word-word relation classification. In The AAAI Conference on Artificial Intelligence, volume 36, pages 10965–10973, Washington, USA. AAAI Press.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In International Conference on Learning Representations, pages 1-19, New Orleans, USA. OpenReview.net.
- Wei Lu and Dan Roth. 2015. Joint mention extraction and classification with mention hypergraphs. In The 2015 Conference on Empirical Methods in Natural Language Processing, pages 857-867, Lisbon, Portugal. Association for Computational Linguistics.
- Da Luo, Yanglei Gan, Rui Hou, Qiao Liu, Tingting Dai, Yuxiang Cai, and Xiaojun Shi. 2024. Unleashing the power of context: Contextual association network with cross-task attention for joint relational extraction. Expert Systems with Applications, 238:121866.
- Xuezhe Ma and Eduard Hovy. 2016. End-to-end sequence labeling via bi-directional LSTM-CNNs-CRF. In The 54th Annual Meeting of the Association for Computational Linguistics, volume 1, pages 1064-1074, Berlin, Germany. Association for Computational Linguistics.
- Aldrian Obaja Muis and Wei Lu. 2017. Labeling gaps between words: Recognizing overlapping mentions with mention separators. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pages 2608–2618, Copenhagen, Denmark. Association for Computational Linguistics.
- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2013. On the difficulty of training recurrent neural networks. In International Conference on Machine Learning, pages 1310–1318, Atlanta, USA. Pmlr, PMLR.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yangi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of Machine Learning Research, 21(1):1-67.

Roberts and G. Lawrence. 1965. Machine Perception of Three-Dimensional Solids. Thesis, MIT: Massachusetts Institute of Technology.

727

728

729

730

731

732

733

734

735

736

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

782

- Lei Sha, Feng Qian, Baobao Chang, and Zhifang Sui. 2018. Jointly extracting event triggers and arguments by dependency-bridge RNN and tensor-based argument interaction. In The Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence, volume 32, pages 5916-5923, Louisiana, USA. AAAI Press.
- Yongliang Shen, Xinyin Ma, Zeqi Tan, Shuai Zhang, Wen Wang, and Weiming Lu. 2021. Locate and label: A two-stage identifier for nested named entity recognition. In The 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, volume 1, pages 2782–2794, Online. Association for Computational Linguistics.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023a. Diffusion-NER: Boundary diffusion for named entity recognition. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3875-3890, Toronto, Canada. Association for Computational Linguistics.
- Yongliang Shen, Zeqi Tan, Shuhui Wu, Wenqi Zhang, Rongsheng Zhang, Yadong Xi, Weiming Lu, and Yueting Zhuang. 2023b. PromptNER: Prompt locating and typing for named entity recognition. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 12492–12507, Toronto, Canada. Association for Computational Linguistics.
- Takashi Shibuya and Eduard Hovy. 2020. Nested named entity recognition via second-best sequence learning and decoding. Transactions of the Association for Computational Linguistics, 8:605-620.
- Mohammad Golam Sohrab and Makoto Miwa. 2018. Deep exhaustive model for nested named entity recognition. In The Conference on Empirical Methods in Natural Language Processing, pages 2843-2849, Brussels, Belgium. Association for Computational Linguistics.
- Jana Straková, Milan Straka, and Jan Hajic. 2019. Neural architectures for nested NER through linearization. In The 57th Annual Meeting of the Association for Computational Linguistics, pages 5326-5331, Florence, Italy. Association for Computational Linguistics.
- Chuanqi Tan, Wei Qiu, Mosha Chen, Rui Wang, and Fei Huang. 2020. Boundary enhanced neural span classification for nested named entity recognition. In The Thirty-Fourth AAAI Conference on Artificial Intelligence, volume 34, pages 9016–9023, New York, USA. AAAI Press.

Zeqi Tan, Yongliang Shen, Shuai Zhang, Weiming Lu,

and Yueting Zhuang. 2021a. A sequence-to-set net-

work for nested named entity recognition. In The

Thirtieth International Joint Conference on Artificial

Intelligence, IJCAI-21, pages 3936–3942, Online.

International Joint Conferences on Artificial Intelli-

Zeqi Tan, Yongliang Shen, Shuai Zhang, Weiming Lu,

and Yueting Zhuang. 2021b. A sequence-to-set net-

work for nested named entity recognition. In Pro-

ceedings of the Thirtieth International Joint Confer-

ence on Artificial Intelligence, IJCAI-21, pages 3936-

3942. International Joint Conferences on Artificial

Wei Tang, Benfeng Xu, Yuyue Zhao, Zhendong Mao, Yifeng Liu, Yong Liao, and Haiyong Xie. 2022.

UniRel: Unified representation and interaction for

joint relational triple extraction. In Proceedings of

the 2022 Conference on Empirical Methods in Nat-

ural Language Processing, pages 7087–7099, Abu

Dhabi, United Arab Emirates. Association for Com-

Juncheng Wan, Dongyu Ru, Weinan Zhang, and Yong

Yu. 2022. Nested named entity recognition with span-

level graphs. In The 60th Annual Meeting of the Association for Computational Linguistics, volume 1,

pages 892-903, dublin, Ireland. Association for Com-

Jue Wang and Wei Lu. 2020. Two are better than

one: Joint entity and relation extraction with tablesequence encoders. In Proceedings of the 2020 Con-

ference on Empirical Methods in Natural Language

Processing (EMNLP), pages 1706–1721, Online. As-

Jue Wang, Lidan Shou, Ke Chen, and Gang Chen. 2020.

Pyramid: A layered model for nested named entity

recognition. In The 58th Annual Meeting of the Association for Computational Linguistics, pages 5918-

5928, Online. Association for Computational Lin-

Yiran Wang, Hiroyuki Shindo, Yuji Matsumoto, and

Taro Watanabe. 2021. Nested named entity recog-

nition via explicitly excluding the influence of the

best path. In The 59th Annual Meeting of the Asso-

ciation for Computational Linguistics and the 11th

International Joint Conference on Natural Language

Processing, pages 3547–3557, Online. Association

Yongxiu Xu, Heyan Huang, Chong Feng, and Yue

Hu. 2021. A supervised multi-head self-attention

network for nested named entity recognition. In

Thirty-Fifth AAAI Conference on Artificial Intelli-

gence, AAAI 2021, Thirty-Third Conference on In-

novative Applications of Artificial Intelligence, IAAI

2021, The Eleventh Symposium on Educational Ad-

for Computational Linguistics.

sociation for Computational Linguistics.

Intelligence Organization. Main Track.

gence Organization.

putational Linguistics.

putational Linguistics.

guistics.

- 791
- 793

- 797
- 802

- 810
- 811
- 812

- 818
- 819

825

- 827
- 831

832

833 834

835

836

839 vances in Artificial Intelligence, EAAI 2021, pages 14185–14193, Online. AAAI Press.

Hang Yan, Tao Gui, Junqi Dai, Qipeng Guo, Zheng Zhang, and Xipeng Qiu. 2021. A unified generative framework for various NER subtasks. In The 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint *Conference on Natural Language Processing*, pages 5808–5822, Online. Association for Computational Linguistics.

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

884

885

886

887

888

889

890

891

893

894

- Hang Yan, Yu Sun, Xiaonan Li, and Xipeng Qiu. 2023. An embarrassingly easy but strong baseline for nested named entity recognition. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 1442-1452, Toronto, Canada. Association for Computational Linguistics.
- Bishan Yang and Tom Mitchell. 2016. Joint extraction of events and entities within a document context. In The 15th Annual Conference of the North American Chapter of the Association for Computational Linguistics, pages 289-299, San Diego, USA. Association for Computational Linguistics.
- Shunyu Yao, Kai Shuang, Rui Li, and Sen Su. 2021. Fgcan: Filter-based gated contextual attention network for event detection. Know.-Based Syst., 228(C).
- Juntao Yu, Bernd Bohnet, and Massimo Poesio. 2020. Named entity recognition as dependency parsing. In The 58th Annual Meeting of the Association for Computational Linguistics, pages 6470-6476, Online. Association for Computational Linguistics.
- Zheng Yuan, Chuangi Tan, Songfang Huang, and Fei Huang. 2022. Fusing heterogeneous factors with triaffine mechanism for nested named entity recognition. In Findings of the Association for Computational Linguistics: ACL 2022, pages 3174-3186, dublin, Ireland. Association for Computational Linguistics.
- Sheng Zhang, Hao Cheng, Jianfeng Gao, and Hoifung Poon. 2023. Optimizing bi-encoder for named entity recognition via contrastive learning. In The Eleventh International Conference on Learning Representations.
- Shuai Zhang, Yongliang Shen, Zeqi Tan, Yiquan Wu, and Weiming Lu. 2022. De-bias for generative extraction in unified NER task. In The 60th Annual Meeting of the Association for Computational Linguistics, pages 808-818, dublin, Ireland. Association for Computational Linguistics.
- Weixiang Zhao, Yanyan Zhao, Xin Lu, and Bing Qin. 2021. An aspect-centralized graph convolutional network for aspect-based sentiment classification. In Natural Language Processing and Chinese Computing: 10th CCF International Conference, NLPCC 2021, Qingdao, China, October 13-17, 2021, Proceedings, Part II, page 260–271, Berlin, Heidelberg. Springer-Verlag.

 Qinghua Zheng, Yuefei Wu, Guangtao Wang, Yanping Chen, Wei Wu, Zai Zhang, Bin Shi, and Bo Dong.
 2023. Exploring interactive and contrastive relations for nested named entity recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:2899–2909.

899

900

901

902

903

905

906

907

908

909

910

911 912

- Zexuan Zhong and Danqi Chen. 2021. A frustratingly easy approach for entity and relation extraction. In 2021 Annual Conference of the North American Chapter of the Association for Computational Linguistics, pages 50–61, Online. Association for Computational Linguistics.
- Enwei Zhu and Jinpeng Li. 2022. Boundary smoothing for named entity recognition. In *The 60th Annual Meeting of the Association for Computational Linguistics*, pages 7096–7108, dublin, Ireland. Association for Computational Linguistics.

Table 6: Statistics of the datasets used in the experiments. The "Len" column represents the average length of sentences or entities in each dataset.

		Train	Dev	Test	Len	Overlap rate
ACE04	Sen	6,297	742	824	23.52	45 (90)
	Ent	22,231	2,514	3,036	2.64	45.08%
ACE05	Sen	7,178	960	1,051	20.59	20 110/-
	Ent	25,300	3,321	3,099	2.40	39.11%
GENIA	Sen	15,023	1,669	1,854	25.41	18 /1%
ULINIA	Ent	45,144	5,365	5,506	1.97	10.4170

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

A Data Statistics

The ACE04 and ACE05 datasets contain seven entity types: Person (PER), Organization (ORG), Geo-Political Entity (GPE), Location (LOC), Facility (FAC), Weapon (WEA), and Vehicle (VEH). The GENIA datasets including five categories: DNA, RNA, Protein, Cell line, and Cell type. According to statistical analysis, 30% of the sentences in the ACE04 and ACE05 datasets contain nested entities, while the GENIA dataset has 17% of sentences with nested entities. The statistical information of the three benchmark datasets is shown in Table 6.

It is worth emphasizing that Yan et al. (2023) observed that despite the usage of the same dataset in recent studies (Wan et al., 2022; Zhu and Li, 2022; Yuan et al., 2022; Li et al., 2022), the statistics of the training datasets differ due to variations in preprocessing methods. Consequently, it would be unfair to directly compare model performance using different versions. In order to address this concern, we utilized the preprocessing code provided by(Yan et al., 2023) and applied it to our dataset. Subsequently, we re-implemented several baseline models in 2022 using the preprocessed dataset and publicly available code. The performance metrics of these models are recorded in Table 1. However, due to the unavailability of code and limited model details, we were unable to fully replicate the Span Graph(Wan et al., 2022) and De-bias(Zhang et al., 2022) models.

B Baseline Details

We compare our method with the following baselines:

1) **Biaffine**: Yu et al. (2020)used a biaffine model to identify nested named entities, predicting the named entity boundaries by predicting the dependency relationship between two words.

2) Second-Best: Wang et al. (2021) recognized

nested entities by explicitly excluding the influence of the optimal path of the probability graph.

3) Seq2Seq: Yan et al. (2021) used a pointerbased approach to convert the entity tagging task into a sequence generation task.

4) Sequence2Set: Tan et al. (2021a) proposed a novel neural network architecture for set prediction specifically for nested NER.

5) De-bias: Zhang et al. (2022) analyzed the incorrect biases in the generation process and used the intra- and inter-entity de-confounding data augmentation methods, to reduce the model's bias.

6) W2NER: Li et al. (2022) modeled unified NER as word-word relationship classification, avoiding conflicts between labels in traditional sequence labeling methods.

7) Locate-and-Label: Shen et al. (2021) modeled the nested NER task as a joint task of entity boundary regression and span classification, improving the training and inference efficiency.

8) BS: Zhu and Li (2022) proposed a boundary smoothing method, which reassigns probabilities from annotated spans to the surrounding ones, to improve the performance of NER models.

9) CNNNER: Yan et al. (2023) used CNN to model the spatial relationships in the score matrix to solve the nested named entity recognition task.

10) Triaffine: Yuan et al. (2022) improved entity recognition performance by obtaining various interaction information between heterogeneous elements such as tokens, entity types, and boundaries.

11) Sequence2Set: Tan et al. (2021a) proposed a novel neural network architecture for set prediction specifically for nested NER.

12) Span-Graph: Wan et al. (2022) modeled nested NER using a span-based graph structure, where each span is represented as a node and spans are connected by edges to enhance the semantic representation capability of the spans.

13) DiffusionNER: Shen et al. (2023a) used the diffusion model for NER task, generating entities by progressive boundary refinement over the noisy spans.

14) **PromptNER**: Shen et al. (2023b) designs a dual-slot multi-prompt template with the position slot and type slot to prompt locating and typing respectively.

15) DINDER: Zhang et al. (2023) frame NER as a representation learning problem that maximizes the similarity between the vector representations of entity mentions and their types.

Table 7: Hyper-parameter settings on different benchmarks

	ACE04	ACE05	GENIA
Batchsize	48	48	8
Epoch	80	80	10
Learning rate	2e-5	2e-5	7e-6
Biaffine size	120	120	400
CNN channel dim	120	120	200
Dropout rate	0.2	0.2	0.1

16) ICR: Zheng et al. (2023) introduces a scale transformation mechanism and a supervised contrastive learning loss to explore interactive and contrastive relations among spans.

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

C Hyper-parameter Details

We utilize RoBERTa-Base (Liu et al., 2019) and BERT-Large (Devlin et al., 2019) as the pre-trained models for the ACE dataset, with a hidden layer size of 768. For the GENIA dataset, we employ BioBERT-Base-v1.1 (Lee et al., 2020) as the pretrained model, also with a hidden layer size of 768. In the SDM module, the number of layers N + 1is set to 2 for all datasets. In the BFM module, the number of convolution blocks n + 1 is set to 2 for all datasets. Except for the extra annotation, the size of Conv used in the model is 3×3 . To minimize memory usage, the SAD operator employed in each layer of the SDM shares parameters, except having different fixed weight templates. The hyper-parameters for the biaffine model were chosen based on the study conducted by (Yan et al., 2023), which also incorporates the multi-head biaffine attention mechanism in its implementation. We set the number of heads to 4 and introduce a span width embedding with a size of 25. By default, the temperature parameter in the Gumbel Softmax estimator is set to 1.

Our model is trained using the AdamW optimizer (Loshchilov and Hutter, 2019). To control overfitting, the L2 norm of the gradient is limited to within 5 by gradient clipping (Pascanu et al., 2013), employed by our model. In the first 10% of the training steps, we gradually increased the learning rate using a linear warm-up scheduler. After the warm-up period, we gradually reduced the learning rate using a linear decay scheduler. All experiments are conducted on NVIDIA Tesla A100 (80G). Other hyper-parameters that vary depending on the datasets are detailed in Table 7.



Figure 3: illustrations of the semantic similarity heatmaps for the entity "*a senior member of the Mukhabarat, Saddam's dreaded secret police*" in Sentence 1. The heatmaps compare two cases: "w/o SDM and BFM" (without Semantic Difference Modeling and Boundary Fusion Module) and "DiFiNet" (with SDM and BFM).



Figure 4: Semantic similarity visualization. It illustrates the semantic similarity between the entity "*them*" and other entities in Sentence 3. Each column has a similar meaning to the corresponding column in Figure 3.

D Semantic Similarity Visualization and Analysis of Cases

1042

1043

1044

1045

1046

1047

1048

1049 1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

In order to visualize the semantic similarity between example instances from Table 5, we provide corresponding semantic similarity heatmaps. Specifically, Figure 3 and Figure 4 display partial semantic similarity heatmaps for instance 1 and instance 3, respectively. To ensure optimal clarity, we present the complete semantic similarity heatmap for instance 2, as shown in Figure 5 and Figure 6. Within these heatmaps, each value within a color block represents the cosine similarity of the span semantic vectors.

Figure 3 demonstrates the effectiveness of appropriately modeling semantic difference information between spans in addressing the issue of boundary insensitivity between nested entities. By utilizing SDM and BFM, the semantic similarity between different nested entities decreases, facilitating their differentiation by the classifier. For instance, in Figure 3, the entity "*a senior* ... *police*" with PER type exhibits a 15% decrease in semantic similarity with the entity "*the Mukhabarat*" of ORG type.

In Figure 4, the explicit incorporation of span semantic difference information is shown to enhance



Figure 5: The similarity heatmap of span semantics generated by DiFiNet without SDM and BFM. It illustrates the semantic similarity between the entity "*fred fred*" and other entities in Sentence 1. The tokens on the vertical axis represent the starting tokens of the spans, while the tokens on the horizontal axis represent the ending tokens of the spans.



Figure 6: Display of the similarity heatmap of span semantics generated by DiFiNet, with the same vertical and horizontal axis settings as Figure 5.

the semantic representation capability of DiFiNet, leading to improved overall robustness. In the absence of SDM and BFM, the similarity between "*them*" and other entities tends to be relatively high, with over half of the entities displaying a similarity of 70% or higher. However, with the inclusion of semantic difference information, the similarity between entities decreases significantly. Even the highest semantic similarity, which occurs with the same type entity "*sudarsono*", remains below 70%.

The visualization results of sentence 2 (Figure 5 and 6) further validate the aforementioned observation. The original model faces difficulties in distinguishing the named entity "*fred fred*" from non-entity spans when confronted with the unseen boundary word "*fred*", leading to a boundary insen-

1082

1067

1083	sitivity problem. In contrast, by leveraging the guid-
1084	ance of span semantic difference information, the
1085	model accurately identifies "fred fred" as a named
1086	entity of type PER and successfully distinguishes
1087	it from the nearly identical span "fred".