# Spans, Not Tokens: A Span-Centric Model for Multi-Span Reading Comprehension

Anonymous ACL submission

#### Abstract

Multi-span reading comprehension (MSRC) requires machines to extract multiple noncontiguous spans from a given context to answer a question. Existing MSRC methods either predict the positions of the start and end tokens of answer spans, or predict the BIO tag of each token. Such token-centric paradigms can hardly capture dependencies among spans which are critical to MSRC. In this paper, we propose a span-centric scheme where spans, as opposed to tokens, are directly represented and scored to qualify as answers. Thanks to the explicit representation of spans in the scheme, our implementation called SpanQualifier beneficially models intra-span and inter-span interactions. Our extensive experiments on three MSRC datasets demonstrate the effectiveness of our span-centric scheme and show that Span-Qualifier achieves state-of-the-art results.

# 1 Introduction

014

017

024

032

Machine reading comprehension (MRC) aims to automatically answer questions from a given context, which is a challenge task of natural language understanding and a basic module of many question answering systems (Chen et al., 2017; Liu et al., 2020). MRC was usually regarded as a singlespan extraction task (Rajpurkar et al., 2016, 2018; Kwiatkowski et al., 2019); single-span MRC models extract only one span as the answer (Seo et al., 2017; Wang and Jiang, 2017; Yu et al., 2018).

**Task:** It may be unrealistic to assume a continuous single span as a good answer to all questions. A more general task is multi-span reading comprehension (MSRC) (Zhu et al., 2020; Ju et al., 2022;

**Question:** Who were the greek philosophers who contributed the basic information about atoms?

**Context:** <u>Democritus</u> was an Ancient Greek pre-Socratic philosopher primarily remembered today for his formulation of an atomic theory of the universe. His exact contributions are difficult to disentangle from those of <u>his mentor Leucippus</u>, as they are often mentioned together in texts. ...

Answer spans: { Democritus, his mentor Leucippus }

Figure 1: An example of MSRC task. Answer spans in the context are underlined.

Cui et al., 2021; Li et al., 2022). As shown in Figure 1, MSRC requires extracting an unspecified number (two in this example) of non-contiguous spans from a given context to answer a question. Moreover, spans are not independently extracted, for example, the second answer span "his mentor Leucippus" is dependent on the first answer span "Democritus", which needs to be captured.

**Challenges:** Existing methods for MSRC either predict the positions of the start and end tokens of answer spans by adapting single-span MRC models to MSRC (Hu et al., 2019; Yang et al., 2021), or treat MSRC as a sequence tagging task by predicting the Beginning-Inside-Outside (BIO) tag of each token (Segal et al., 2020; Li et al., 2022). We regard these methods as *token-centric* which define features, train models, and predict answers from the perspective of tokens. Their effectiveness is limited in two aspects. Firstly, it is non-trivial and often suboptimal to convert token-level predictions to multiple spans as answers. Secondly, token-level 035

036

representation and optimization are not conducive to capturing dependencies among spans.

057

061

062

067

074

077

079

086

090

100

101

102

103

104

**Our Work:** To overcome the limitations, we propose a *span-centric* scheme for MSRC and present an effective implementation called SpanQualifier. In response to the first limitation, our span-centric scheme generates representations for all spans in the context and predicts a qualification threshold. Spans are scored, and a span will qualify as an answer if its score exceeds the qualification threshold, thus avoiding the conversion from token-level predictions. For the second limitation, SpanQualifier interacts intra-span and inter-span representations to capture their dependencies and, moreover, a global loss function is designed to jointly optimize over all spans and predict their scores.

To summarize, our contributions include

- a novel span-centric scheme for MSRC which more naturally fits the task and empirically outperforms existing token-centric methods,
- an implementation of the scheme which models intra/inter-span interactions and obtains state-of-the-art results on three datasets.

Our code and data are in the supplementary material and will be open on GitHub after acceptance.

### 2 Related Work

### 2.1 Pointer Paradigm for MSRC

The Pointer paradigm based on Pointer Network (Vinyals et al., 2015) is one of the most popular schemes for single-span MRC, first proposed by Wang and Jiang (2017), and various successors extend it with sophisticated matching mechanisms (Seo et al., 2017; Yu et al., 2018). As shown in Figure 2a, the predicted probabilities of being the start or end of an answer span are normalized over all tokens in the context, and a post-processing decoding strategy such as beam search is adopted to select a span as the answer. It is non-trivial to adapt this paradigm to MSRC due to the unspecified number of target spans. MTMSN (Hu et al., 2019) predicts the number of answer spans and extracts non-overlapping spans based on an algorithm in computer vision (Rosenfeld and Thurston, 1971). MUSST (Yang et al., 2021) iteratively extracts and masks a span with the highest probability until meeting a stop symbol. Essentially, these methods predict token-level probabilities and convert them into multiple spans by non-differentiable decoding strategies, which may under-utilize the training





Figure 2: Two existing paradigms for MSRC.

data for MSRC. By contrast, RASOR (Lee et al., 2016) employs span representations to select a single span, which is extended by Chen et al. (2020) with span merging to extract at most two spans. This is still not enough for MSRC where answers are often more than two spans.

105

106

107

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

Unlike the above methods, in our proposed spancentric scheme, all possible spans are scored and may qualify as answers, and our end-to-end model is differentiable to exploit training data.

### 2.2 Tagger Paradigm for MSRC

The Tagger paradigm based on sequence tagging has been applied to many tasks, such as part-ofspeech tagging and information extraction, and first applied to answer extraction by Yao et al. (2013). As shown in Figure 2b, a BIO tag is predicted for each token. This paradigm can be relatively easily adapted to MSRC (Yoon et al., 2022; Li et al., 2022). Further, TASE (Segal et al., 2020) employs a multi-head architecture to separately handle single-span and multi-span questions. Li et al. (2022) incorporates answer structure and answer number prediction into a multi-task model.

Tagger is trained by local normalization on each individual token, which suffers from label bias (Andor et al., 2016). It is less informative than global normalization performed in our implementation over all possible spans to consider their dependencies when calculating our loss function.

### 2.3 Information Extraction

Span-centric models have shown their effectiveness in information extraction tasks, including named entity recognition (NER), relation extraction, and



Figure 3: Overview of SpanQualifier. Dashed boxes represent extensions of the vanilla implementation.

event extraction (Luan et al., 2018, 2019; Wadden et al., 2019). They are trained by calculating a cross-entropy loss over all possible spans.
However, this training objective is not suitable for MSRC due to the extremely unbalanced labels in its training data. Indeed, the context in a MSRC task typically contains tens of thousands of spans while only a few are answers.

Recently, Su et al. (2022) propose a NER model 146 with a universal loss function that extends the soft-147 max and cross-entropy loss for handling label im-148 balance, but their span representations obtained 149 by interacting only between its start and end token representations are not sufficient for MSRC. 151 As pointed out by Ju et al. (2022), different from 152 NER which requires extracting all the spans of 153 each given type, MSRC is more difficult and re-154 lies on a thorough understanding of the context 155 because some spans of a given type are not cor-156 rect answers. To enhance span representation, Tan et al. (2020) add a boundary detection task to predict the start and end tokens of entities, and some 159 works use computationally expensive bi-affine (Yu 160 et al., 2020) or even tri-affine (Yuan et al., 2022) 161 mechanisms. By contrast, our SpanQualifier inexpensively obtains context-dependent span representations by two successive interaction layers where 164 the former uses a supervised attention mechanism 165 to highlight intra-span tokens to be aggregated and 166 the latter captures inter-span dependencies. 167

# 3 Approach

Given a question Q and a context  $C = c_0, \ldots, c_n$ with n tokens, the goal of MSRC is to extract a set of m answer indexes  $A = \{[p_1^s, p_1^e], \ldots, [p_m^s, p_m^e]\}$ from the context, where  $p_i^s$  and  $p_i^e$  are the start and end positions of the *i*-th answer span. 169

170

171

172

173

174

175

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

Figure 3 shows an overview of our SpanQualifier, which simultaneously extracts multiple spans as answers by implementing a span-centric scheme with four stacked modules (Section 3.1), two extra interaction layers within and between spans (Section 3.2), and a loss function combining global normalization and maximum margin (Section 3.3).

#### 3.1 (Vanilla) Span-Centric Scheme

The proposed scheme successively constructs token representations, boundary (i.e., start/end token) representations, and span representations, and then selects answers by predicting a qualification threshold for span score. In the remainder of this section,  $MLP(\cdot)$  represents a two-layer perceptron.

**Token Representation** We firstly capture the dependencies between the context C and the question Q to generate token representations. We feed Q and C into a pre-trained language model (PLM) such as BERT (Devlin et al., 2019) or ALBERT (Lan et al., 2020) in a standard way:

$$\mathbf{H} = \mathsf{PLM}([\mathsf{CLS}] \ Q \ [\mathsf{SEP}] \ C \ [\mathsf{SEP}]), \tag{1}$$

246

247

248

249

251

252

253

255

256

258

260

262

265

266

267

268

269

270

271

272

273

274

275

237

where **H** represents the hidden states of the last layer of PLM. From **H** we extract the global representation  $\mathbf{H}^{\text{CLS}} \in \mathbb{R}^{d_1}$  and the context representation  $\mathbf{H}^{\text{C}} \in \mathbb{R}^{n \times d_1}$ , where  $d_1$  denotes the dimension of representation and n is the length of C.

**Boundary Representation** Then we aggregate the information about each possible span to its two boundary tokens. In our vanilla implementation, we take the context representation  $\mathbf{H}^{C}$  as input and use a boundary enumeration layer to obtain the start and end token representations:

201

203

204

234

236

$$\begin{split} \mathbf{B}^{s}, \mathbf{B}^{e} &= \texttt{BoundaryEnum}(\mathbf{H}^{C}) \\ \text{where} \quad \mathbf{B}^{s} &= \texttt{MLP}(\mathbf{H}^{C}) \text{ and } \mathbf{B}^{e} &= \texttt{MLP}(\mathbf{H}^{C}) \,, \end{split}$$

where  $\mathbf{B}^{s}, \mathbf{B}^{e} \in \mathbb{R}^{n \times d_{1}}$ . Note that an extra intraspan interaction layer can be inserted here to achieve more intensive information aggregation, which will be introduced in Section 3.2.

211Span RepresentationNext, we generate the rep-212resentations for all possible spans. In our vanilla213implementation, the start and end token representa-214tions  $\mathbf{B}^{s}$ ,  $\mathbf{B}^{e}$  are fused by a span enumeration layer215to obtain the representations of all possible spans:

$$\mathbf{M} = \operatorname{SpanEnum}(\mathbf{B}^{s}, \mathbf{B}^{e})$$

$$= \operatorname{LayerNorm}(\mathbf{N}) \qquad (3)$$
where  $\mathbf{N}_{i,j} = \mathbf{B}_{i}^{s} W_{1}^{s} + \mathbf{B}_{j}^{e} W_{1}^{e} + \mathbf{E}_{|j-i|}$ 

where  $\mathbf{M} \in \mathbb{R}^{n \times n \times d_2}$ ,  $\mathbf{M}_{i,j}$  denotes the represen-217 tation of the span starting at the *i*-th token and end-218 ing at the *j*-th token, and  $\mathbf{E}_{|j-i|} \in \mathbb{R}^{d_2}$  is a learn-219 able embedding of span length. LayerNorm $(\cdot)$  is a layer normalization operator (Ba et al., 2016), and  $W_1^s, W_1^e \in \mathbb{R}^{d_1 \times d_2}$  are learnable parameters used to reduce the dimension from  $d_1$  to  $d_2$  to improve 223 efficiency. We empirically set  $d_2 = 64$ . Note that an extra interaction layer for capturing dependencies among spans can be inserted here, which will be introduced in Section 3.2. 227

Span Selection Finally, we score all possible
spans and select answer spans by predicting a qualification threshold. We feed both the span representations M and the global representation H<sup>CLS</sup> into
a scoring layer:

$$\mathbf{S}, qs^{\text{ext}} = \text{Scoring}(\mathbf{M}, \mathbf{H}^{\text{CLS}})$$
  
where  $\mathbf{S} = \text{MLP}(\mathbf{M})$  and  $qs^{\text{ext}} = \text{MLP}(\mathbf{H}^{\text{CLS}})$ , (4)

where  $\mathbf{S} \in \mathbb{R}^{n \times n}$  contains the scores of all possible spans and  $\mathbf{S}_{i,j} \in \mathbf{S}$  denotes the score of the span starting at the *i*-th token and ending at the *j*-th token. The lower triangular portion of S and the portion representing spans containing more than k tokens are masked, where k is the longest length of an answer in the training set.

At inference time, all the spans with score higher than  $qs^{\text{ext}}$  will qualify as answers, where  $qs^{\text{ext}} \in \mathbb{R}$ represents a qualification threshold. If two spans overlap, we will keep the higher scored one.

#### 3.2 Span Interaction

To extend the vanilla implementation, two interaction layers are added to enhance span representations and capture dependencies among spans.

**Intra-Span Interaction** Not only the start and end tokens but also those inside a span are crucial to its representation. Therefore, we insert an intra-span interaction layer into the boundary representation module after Equation (2) to incorporate the attention weights of internal tokens:

$$\hat{\mathbf{B}}^{s}, \hat{\mathbf{B}}^{e} = \text{BoundaryEnum}(\mathbf{H}^{C}), \\ \mathbf{B}^{s}, \mathbf{B}^{e} = \text{IntraSpan}(\hat{\mathbf{B}}^{s}, \hat{\mathbf{B}}^{e}).$$
 (5)

As for IntraSpan $(\cdot, \cdot)$ , our implementation resembles Equation (3) and Equation (4). Start token representations are intensified as follows:

$$\mathbf{M}^{s} = \operatorname{SpanEnum}(\hat{\mathbf{B}}^{s}, \hat{\mathbf{B}}^{e}),$$
  

$$\mathbf{G}^{s}, qs^{s} = \operatorname{Scoring}(\mathbf{M}^{s}, \mathbf{H}^{\operatorname{CLS}}),$$
  

$$\mathbf{B}^{s} = \operatorname{Softmax}(\mathbf{G}^{s})(\hat{\mathbf{B}}^{s}W_{2}^{s}),$$
  
(6)

and end token representations are intensified by

$$\begin{split} \mathbf{M}^{e} &= \mathtt{SpanEnum}(\hat{\mathbf{B}}^{s}, \hat{\mathbf{B}}^{e}), \\ \mathbf{G}^{e}, qs^{e} &= \mathtt{Scoring}(\mathbf{M}^{e}, \mathbf{H}^{\mathrm{CLS}}), \\ \mathbf{B}^{e} &= \mathtt{Softmax}((\mathbf{G}^{e})^{\mathsf{T}})(\hat{\mathbf{B}}^{e}W_{2}^{e}), \end{split}$$
(7)

where  $\mathbf{G}^{s}, \mathbf{G}^{e}$  represent attention weights, and  $W_{2}^{s}, W_{2}^{e} \in \mathbb{R}^{d_{1} \times d_{1}}$  are learnable parameters.

**Inter-Span Interaction** To capture the dependencies among spans which are useful in MSRC, we insert an inter-span interaction layer into the span representation module after Equation (3):

$$\begin{split} \hat{\mathbf{M}} &= \mathtt{SpanEnum}(\mathbf{B}^{s},\mathbf{B}^{e})\,, \\ \mathbf{M} &= \mathtt{InterSpan}(\hat{\mathbf{M}})\,, \end{split} \tag{8}$$

where we implement  $InterSpan(\cdot)$  by an onelayer 2D convolutional neural network (CNN).

#### 3.3 Training

Since answer spans in MSRC correlate with each other, a global qualification loss function is introduced for supervising span selection. It is also adapted to supervise intra-span interaction.

	# Train	# Dev	# Test	Distributi	on by answ	er number/	Average	Average
	# 11aiii	# Dev	# 1051	0	1	$\geq 2$	answer number	context length
MultiSpanQA	4,577	653	653	-	-	100.00%	2.89	251
MultiSpanQA-E	13,731	1,959	1,959	33.43%	32.59%	33.98%	1.30	226
QUOREF	16,083	2,186	2,276	-	90.56%	9.44%	1.14	324

Table 1: Dataset statistics.

Global Qualification Loss The Global Qualification (GQ) loss consists of two parts: the Global Normalization (GN) loss and the Maximum Margin (MM) loss. Considering the correlation between positive spans, we use the GN loss to jointly optimize the scores of all positive spans. Moreover, we use the MM loss to maximize the margin between positive and negative spans.

284

289

290

294

296

Given a set of scores of positive spans  $\Omega^+ \subseteq \mathbf{S}$ , a set of scores of negative spans  $\Omega^- \subseteq \mathbf{S}$ , and a qualification threshold  $qs^{\text{ext}}$ , the GN loss is

$$L_{\rm GN}(\Omega^+, \Omega^-) = -\log \frac{\sum\limits_{s \in \Omega^+} \exp(s)}{\sum\limits_{s \in \Omega^+} \exp(s) + \sum\limits_{s \in \Omega^-} \exp(s)} \cdot \frac{1}{(9)}$$

The objective is to maximize the sum of probabilities of all positive spans. The MM loss comprises

$$\begin{split} \mathbf{L}_{\rm MM}^+(\Omega^+, qs^{\rm ext}) &= {\rm Max}(1-({\rm Min}(\Omega^+)-qs^{\rm ext}), 0)\,, \\ \mathbf{L}_{\rm MM}^-(\Omega^-, qs^{\rm ext}) &= {\rm Max}(1-(qs^{\rm ext}-{\rm Max}(\Omega^-)), 0)\,. \end{split} \eqno(10)$$

The objective is to maximize the soft margin between positive and negative spans. The GQ loss combines the GN loss and the MM loss:

**Supervised Extraction** We employ the GQ loss as our main loss for supervising span extraction:

$$\mathcal{L}_{\text{ext}} = L_{\text{GQ}}(\Omega^+, \Omega^-, qs^{\text{ext}}).$$
(12)

298Supervised AttentionWe also adapt the GQ loss299to supervise the attention weights in the intra-span300interaction layer. Let  $\Psi^{s+} \subseteq \mathbf{G}^s$  and  $\Psi^{e+} \subseteq \mathbf{G}^e$  be301the sets of attention weights of tokens in positive302spans. For each positive span  $[p^s, p^e], \Psi^{s+}$  contains303 $\mathbf{G}_{p^s,p^s}^s, \mathbf{G}_{p^s,p^{s+1}}^s, \dots, \mathbf{G}_{p^s,p^e}^s \in \mathbf{G}^s$ , and  $\Psi^{e+}$  con-304tains  $\mathbf{G}_{p^s,p^e}^s, \mathbf{G}_{p^s+1,p^e}^s, \dots, \mathbf{G}_{p^e,p^e}^s \in \mathbf{G}^e$ . Let  $\Psi^{s-}$ 305and  $\Psi^{e-}$  be the sets of the remaining attention306weights in  $\mathbf{G}^s$  and  $\mathbf{G}^e$ , respectively. Given the307thresholds  $qs^s$  and  $qs^e$  computed by Equation (6)308and Equation (7), respectively, our attention loss is

$$\mathcal{L}_{att} = L_{GQ}(\Psi^{s+}, \Psi^{s-}, qs^{s}) + L_{GQ}(\Psi^{e+}, \Psi^{e-}, qs^{e}).$$
(13)

**Total Loss** Finally, we sum the extraction loss and the attention loss for joint training:

$$\mathcal{L} = \mathcal{L}_{\text{ext}} + \mathcal{L}_{\text{att}} \,. \tag{14}$$

310

311

312

313

314

315

316

317

318

319

321

322

323

324

326

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

# 4 Experimental Setup

#### 4.1 Datasets

We used three MSRC datasets. **MultiSpanQA** and **MultiSpanQA-E**xpanded (Li et al., 2022) are the latest datasets for MSRC created from a question answering dataset (Kwiatkowski et al., 2019). MultiSpanQA only contains multi-span questions, while MultiSpanQA-E also contains unanswerable and single-span questions. **QUOREF** (Dasigi et al., 2019) containing questions requiring coreferential reasoning has been widely used in MSRC research.

Since the official test sets of these datasets are not public, we took their official dev set as our test set and randomly held out an in-house dev set from the official training set. Some statistics about the datasets are shown in Table 1.

### 4.2 Participating Methods

We experimented with two variants of our method: **SpanQualifier** representing the full version and **SpanQualifier**<sub>vanilla</sub> representing the vanilla version without span interaction in Section 3.2.

We compared with two pointer-based methods for MSRC. **MTMSN** (Hu et al., 2019) predicts the number of answer spans and repeatedly extracts non-overlapping spans until reaching the predicted number. **MUSST** (Yang et al., 2021) iteratively extracts spans until meeting a special stop symbol.

We compared with four tagger-based methods for MSRC. **Tagger**vanilla (Li et al., 2022) uses a standard BIO tagger. **Tagger**multi-task (Li et al., 2022) enhances Taggervanilla with answer structure and answer number prediction. **Tagger**crf extends Taggervanilla with a CRF layer (Lafferty et al., 2001) to capture tag dependencies. **TASE** (Segal et al., 2020) employs a multi-head mechanism to separately handle single-span and multi-span questions. We further extended it to separately handle unanswerable questions in MultiSpanQA-E.

	MultiSpanQA			MultiSpanQA-E				QUOREF				
	E	М	Р	М	E	М	Р	М	E	М	Р	М
	dev	test	dev	test	dev	test	dev	test	dev	test	dev	test
BERT												
MTMSN	46.67•▲	43.55•▲	72.17•▲	70.44•▲	49.96•▲	50.60•▲	70.61•▲	70.92	62.26°▲	62.23 🔺	73.81 🔺	73.62
MUSST	51.35**	48.83•▲	74.61•▲	72.36•▲	54.26•▲	54.65•▲	70.81•▲	70.60•▲	33.35•▲	31.91•▲	55.21•▲	54.23•▲
Taggervanilla	60.36•▲	57.89•▲	73.76•▲	72.60•▲	60.39•▲	59.81•▲	72.38 🗅	72.33 🔺	56.33•▲	55.67•▲	68.32•▲	67.47•▲
Tagger <sub>multi-task</sub>	51.63•▲	51.72•▲	67.97•▲	69.06•▲	54.93•▲	56.53•▲	68.54•▲	70.36•▲	39.55•▲	39.00•▲	51.30•▲	49.75•▲
Tagger <sub>crf</sub>	61.28•▲	58.94•▲	74.55•▲	73.43•▲	57.00•▲	57.23•▲	66.94•▲	67.18•▲	56.73•▲	55.09•▲	68.24•▲	66.89•▲
TASE	57.41•▲	54.99•▲	73.17•▲	72.59•▲	51.01•▲	52.47•▲	65.38•▲	66.99•▲	58.58•▲	56.98•▲	69.28•▲	67.50•▲
SpanQualifiervanilla	66.38	64.92	78.61	77.69	63.25	64.24	73.19	73.80	63.92	62.31	73.33	71.69
SpanQualifier	69.94	68.41	81.26	80.61	65.47	66.08	74.09	75.11	67.19	65.65	75.70	74.11
ALBERT												
MTMSN	40.85•▲	38.12•▲	<b>63.17</b> •▲	61.14•▲	41.74•▲	43.90•▲	56.28•▲	58.58•▲	45.78•▲	44.35•▲	57.91•▲	56.64•▲
MUSST	54.13•▲	51.69•▲	77.45•▲	75.56•▲	57.27•▲	57.34•▲	73.95•▲	74.52•▲	41.73•▲	39.27•▲	65.13•▲	63.48•▲
Taggervanilla	67.06•▲	64.78•▲	80.07•▲	79.32•▲	61.73•▲	61.36•▲	72.77•▲	72.73•▲	70.78•▲	<b>68.97</b> •▲	79.68•▲	77.68•▲
Tagger <sub>multi-task</sub>	49.78•▲	48.75•▲	67.99•▲	67.34•▲	51.65•▲	52.15•▲	65.89•▲	66.38•▲	50.95•▲	48.91•▲	61.77•▲	59.04•▲
Tagger <sub>crf</sub>	68.11•▲	65.43•▲	81.03•▲	80.00•▲	61.99•▲	61.75•▲	72.46•▲	72.86•▲	70.63•▲	69.03•▲	79.35•▲	77.75•▲
TASE	65.55•▲	62.06•▲	80.85•▲	79.34•▲	57.52•▲	56.95•▲	70.95•▲	71.11•▲	71.90•▲	70.31•▲	81.39•▲	79.90•▲
SpanQualifiervanilla	68.92	67.23	82.42	81.76	65.68	66.07	75.85	76.53	73.90	71.77	81.95	79.92
SpanQualifier	72.29	70.20	83.62	82.68	67.24	66.99	76.79	77.01	75.18	72.84	82.56	80.68

Table 2: Comparison with baselines. We mark the results of baselines that are significantly lower than the results of SpanQualifier<sub>vanilla</sub> under p < 0.01 (•) or p < 0.05 (°), and of SpanQualifier under p < 0.01 (•) or p < 0.05 (<sup> $\Delta$ </sup>).

# 4.3 Evaluation Metrics

351

357

358

361

367

370

372

373

377

378

We followed Li et al. (2022) to use two metrics. Exact Match F1 (**EM**) measures the percentage of predicted spans that exactly match any gold-standard answer span. Partial Match F1 (**PM**) generalizes EM by using the length of the longest common substring to score each predicted span based on its nearest gold-standard answer span.

#### 4.4 Implementation Details

We used two PLMs of different sizes: BERT-baseuncased (Devlin et al., 2019) and ALBERT-basev2 (Lan et al., 2020). We implemented our models with PyTorch 1.10 and HuggingFace Transformers 4.12.5. The dimension of the hidden layer of MLP was 768. We implemented CNN using kernel size =  $5 \times 5$ , stride = 1, and in-channels = out-channels = 64. We used the Adam optimizer and set warmup fraction = 0.1, weight decay = 0.01, maximum length = 512, and epoch = 30. For each model, we searched for the best learning rate from  $\{1e - 5, 3e - 5\}$ , and batch size from  $\{24, 32\}$ . We used three seeds  $\{0, 1, 2\}$  and took the mean results. All the experiments were performed on RTX 3090 (24G).

# **5** Experimental Results

### 5.1 Comparison with Baselines

In Table 2, SpanQualifer achieves the best results in all the settings. It significantly outperforms all the baselines under p < 0.01 except in one setting under p < 0.05. Satisfyingly, SpanQualifier<sub>vanilla</sub>

	MultiSp	MultiSpanQA-E		REF
	EM	PM	EM	PM
BERT				
MTMSN	35.86	64.85	57.69	75.34
MUSST	47.20	67.23	29.79	57.15
Taggervanilla	59.61	73.52	61.66	74.43
Tagger <sub>multi-task</sub>	54.64	70.58	42.11	56.02
Tagger <sub>crf</sub>	56.99	68.51	60.85	73.81
TASE	49.41	65.43	61.47	74.17
SpanQualifiervanilla	63.53	74.35	64.57	74.75
SpanQualifier	65.59	75.98	68.82	77.98
ALBERT				
MTMSN	27.58	47.48	36.48	54.20
MUSST	50.40	71.60	33.12	64.17
Taggervanilla	61.58	74.92	71.44	82.55
Tagger <sub>multi-task</sub>	48.23	65.59	49.99	65.16
Tagger <sub>crf</sub>	62.11	75.03	71.70	83.15
TASE	53.37	69.92	70.55	81.80
SpanQualifiervanilla	64.67	76.29	72.38	81.67
SpanQualifier	65.81	76.91	72.31	81.83

Table 3: Comparison with baselines on multi-span questions in the test sets of MultiSpanQA-E and QUOREF.

also exceeds most baselines, demonstrating the effectiveness of our vanilla span-centric scheme.

By excluding single-span and unanswerable questions from MultiSpanQA-E and QUOREF, in Table 3, the gap between SpanQualifier and pointerbased methods becomes larger on multi-span questions. Indeed, by contrast with pointer-based methods, SpanQualifier and tagger-based methods achieve better results on multi-span questions than on other questions, especially on QUOREF, indicating their better suitability for the MSRC task.

# 5.2 Ablation Study: Span Representation

To analyze the effectiveness of our span-centric scheme, we implemented TokenQualifier which is

381

382

	MultiS	panQA	MultiS	panQA-E	QUC	DREF
	EM	PM	EM	PM	EM	PM
BERT						
SpanQualifiervanilla	64.92	77.69	64.24	73.80	62.31	71.69
TokenQualifier	61.35	75.52	59.81	72.33	60.79	70.25
ALBERT						
SpanQualifier <sub>vanilla</sub>	67.23	81.76	66.07	76.53	71.77	79.92
TokenQualifier	66.28	80.88	63.52	75.67	70.58	78.89

Table 4: Ablation study of span representation on the test sets of three datasets.

	MultiSpanQA		MultiSp	anQA-E	QUOREF	
	EM	PM	EM	PM	EM	PM
BERT						
SpanQualifier	68.41	80.61	66.08	75.11	65.65	74.11
w/o intra-span	68.20	79.84	65.00	73.69	65.49	74.01
w/o inter-span	64.76	77.65	65.61	75.10	62.94	72.27
w/o $\mathcal{L}_{att}$	67.15	78.96	64.12	72.83	65.09	73.66
standard attention	67.89	79.86	65.88	75.16	65.15	73.97
ALBERT						
SpanQualifier	70.20	82.68	66.99	77.01	72.84	80.68
w/o intra-span	69.14	81.73	67.11	76.81	72.60	80.30
w/o inter-span	68.28	81.88	65.30	76.14	71.51	79.74
w/o $\mathcal{L}_{att}$	68.87	81.56	66.08	75.80	71.89	80.05
Standard Attention	68.12	81.13	65.75	76.72	72.12	80.01

Table 5: Ablation study of span interaction on the test sets of three datasets.

a variant of SpanQualifier<sub>vanilla</sub> by removing boundary and span representations and directly using token representations for prediction. In Table 4, EM and PM drop by 0.95–4.43 and 0.86–2.17, respectively, demonstrating the usefulness of our boundary and span representations.

### 5.3 Ablation Study: Span Interaction

396

397

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

Our span interaction consists of an intra-span layer with a supervised attention mechanism and an interspan layer. To analyze their effectiveness, we conducted an ablation study by separately removing each component. We also replaced our attention mechanism for computing  $\mathbf{G}^{s}$  and  $\mathbf{G}^{e}$  in Equation (6) and Equation (7), respectively, with a standard attention mechanism where the key and value are MLP( $\hat{\mathbf{B}}^{s}$ ) and MLP( $\hat{\mathbf{B}}^{e}$ ), respectively.

In Table 5, the full version of SpanQualifier outperforms its variants in almost all the settings, demonstrating the usefulness of the removed/replaced components. In particular, on MultiSpanQA and QUOREF, inter-span interaction appears more influential, while on MultiSpanQA-E, attention supervision is more helpful.

### 5.4 Official Leaderboards

We submitted the results predicted by Span-Qualifier on each official test set to the leaderboard: to MultiSpanQA and MultiSpanQA-E<sup>1</sup> on

	MultiS	MultiSpanQA		panQA-E
	EM	PM	EM	PM
SpanQualifier (RoBERTa Large)	72.75	85.41	70.85	81.64
Taggervanilla (RoBERTa Large)	69.19	83.86	69.05	78.97
Tagger+LIQUID (RoBERTa Base)	67.40	81.16	-	-
SpanQualifier (BERT Base)	64.87	78.79	64.39	74.22
Tagger <sub>multi-task</sub> (BERT Base)	59.28	76.26	42.26	70.47
Taggervanilla (BERT Base)	56.45	75.22	41.38	70.10
Single-span (BERT Base)	14.41	67.56	12.26	67.73

Table 6: Official leaderboards of MultiSpanQA and MultiSpanQA-E.

	QUO	REF
	Global EM	Global F1
SpanQualifier (CorefRoBERTa Large)	81.36	87.24
SpanQualifier (RoBERTa Large)	81.24	87.05
TASE (CorefRoBERTa Large)	80.61	86.70
TASE <sub>CoNLL-joint-qgen</sub>	80.45	86.46
TASE (RoBERTa Large)	79.66	86.13
CorefAdditive	79.11	85.86
CorefRoBERTa Large	75.80	82.81

Table 7: Official leaderboard of QUOREF.

10/11/2022, and to QUOREF<sup>2</sup> on 10/12/2022. For a fair comparison, we used the same PLMs as other methods on the leaderboard: BERT Base (Devlin et al., 2019) and RoBERTa Large (Liu et al., 2019) on MultiSpanQA and MultiSpanQA-E, and RoBERTa Large and CorefRoBERTa Large (Ye et al., 2020) on QUOREF.

Table 6 compares SpanQualifier with other topranked single-model methods on the leaderboards of MultiSpanQA and MultiSpanQA-E at the time of our submission. SpanQualifier outperforms all the methods using the same PLM, including **Tagger+LIQUID** (Lee et al., 2023) whose implementation is not public, and **Single-span** (Li et al., 2022) which is a pointer-based method selecting multiple spans by tuning a threshold on the dev set.

Table 7 compares SpanQualifier with other topranked single-model methods on the leaderboard of QUOREF at the time of our submission. Note that this leaderboard uses Global Exact Match (Global EM) (Rajpurkar et al., 2016) and Global F1 (Dua et al., 2019) as evaluation metrics. Span-Qualifier outperforms all the methods, including TASE<sub>CoNLL-joint-qgen</sub> which jointly trains TASE on QUOREF and CoNLL (Sang and Meulder, 2003), CorefAdditive (Zhang and Zhao, 2022) and CorefRoBERTa Large (Ye et al., 2020) which are pointer-based methods predicting the number of answer spans as in MTMSN.

<sup>&</sup>lt;sup>1</sup>https://multi-span.github.io/

<sup>&</sup>lt;sup>2</sup>https://leaderboard.allenai.org/ quoref/submissions/public

Question	Context	Prediction
Which two hormones	His theory states that preparing the animal for fighting or fleeindull. More	{norepinephrine,
are known for their	specifically, the adrenal medulla produces a hormonal cascade that results in	epinephrine, estro-
roles in the bluefight-	the secretion of catecholamines, especially norepinephrine and epinephrine.	gen, testosterone,
or-flight response?	The hormones estrogen, testosterone, and cortisol, as well as also affect	cortisol}
	how organisms react to stress	
Where does the red	Red hair occurs more frequently in people of northern or western European	{northern or western
hair gene come from?	ancestry, and less frequently in other populations. Red hair appears most	European ancestry,
	commonly in people with two copies of a recessive allele on chromosome 16	other populations}
	which produces an altered version of the MC1R protein	

Table 8: Answer spans predicted by SpanQualifier on MultiSpanQA. Gold-standard answer spans are underlined.



Figure 4: Top spans scored by SpanQualifier on Multi-SpanQA. Gold-standard answer spans are underlined.

#### 5.5 Case Study

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

Figure 4 illustrates the top-scored spans predicted by SpanQualifier for a question in MultiSpanQA. The top-two spans with scores higher than the predicted qualification threshold exactly match the gold standard. By contrast, the pointer-based MUSST predicted a long span containing 52 words, and Tagger<sub>crf</sub> predicted three single-word spans. While their spans partially match the gold standard, their inexactness reveals the limitation of tokencentric methods in MSRC, particularly the suboptimality of converting token-level predictions to multiple answer spans, which has been addressed by our span-centric scheme.

#### 5.6 Error Analysis

We manually analyzed errors made by SpanQualifier on MultiSpanQA and identified two common types of errors in MSRC: incorrect homogeneous answers and missing heterogeneous answers.

For the first type of error, as illustrated in the upper part of Table 8, the gold-standard answer spans are surrounded by many other spans of the same type, which misled our model to pay attention

	BERT	ALBERT
MTMSN	115M	17M
MUSST	109M	12M
Taggervanilla	109M	12M
Tagger <sub>multi-task</sub>	117M	20M
Tagger <sub>crf</sub>	109M	12M
TASE	111M	13M
SpanQualifiervanilla	112M	15M
SpanQualifier	115M	18M

Table 9: Model size (number of parameters).

to their dependencies and extract incorrect answers such as "estrogen", "testosterone", and "cortisol". 474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

503

For the second type of error, as illustrated in the lower part of Table 8, the question has multiple heterogeneous answers of distinct types, including population, chromosome, and protein. Their dependencies are difficult to be captured so that our model failed to find some of them.

# 5.7 Model Size and Inference Time

Table 9 compares model sizes. Despite its superior performance, SpanQualifier has a number of parameters comparable to those of the baselines.

BERT-based SpanQualifier<sub>vanilla</sub> and SpanQualifier use an average of 29.69ms and 33.72ms to answer a question, respectively. ALBERT-based SpanQualifier<sub>vanilla</sub> and SpanQualifier use an average of 36.87ms and 42.15ms, respectively. The inference time of our models is reasonably fast.

### 6 Conclusion

In this paper, we explore the effectiveness of spancentric scheme for MSRC and propose a novel extraction model SpanQualifier. Experiments show that our span-centric scheme outperforms existing pointer-based and tagger-based methods which follow a token-centric scheme and, furthermore, with span interaction our model achieves state-of-the-art results on three MSRC datasets. In the future, we plan to add more expressive interaction methods to the span representation module and to soft answer labels to accommodate overlapping spans.

# 561 563 564 565 566 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586 587 588 589 591 592 594 595 596 597 598 599 600 601 602 603 604 605 606 607 608 609 610

611

612

613

614

615

557

558

559

560

# 504 Limitations

505Our model was specifically designed for MSRC.506Some modules (e.g., inter-span interaction) may507not be useful for conventional single-span reading508comprehension. Moreover, considering the poten-509tially different data distributions between single-510span and multi-span questions, our model trained511on a hybrid of single-span and multi-span ques-512tions (e.g., on QUOREF in our experiments) may513not outperform dedicatedly trained models.

#### References

514

515

516

517

518

519

520

521

522

523 524

525

526

527

528

529

530

532

533

534

535

538

539

540

541

542

543

545

546

547

548

549

550

551

552

554

555

- Daniel Andor, Chris Alberti, David Weiss, Aliaksei Severyn, Alessandro Presta, Kuzman Ganchev, Slav Petrov, and Michael Collins. 2016. Globally normalized transition-based neural networks. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers. The Association for Computer Linguistics.
- Lei Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. Layer normalization. *CoRR*, abs/1607.06450.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading wikipedia to answer opendomain questions. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 -August 4, Volume 1: Long Papers, pages 1870–1879. Association for Computational Linguistics.
- Xinyun Chen, Chen Liang, Adams Wei Yu, Denny Zhou, Dawn Song, and Quoc V. Le. 2020. Neural symbolic reader: Scalable integration of distributed and symbolic representations for reading comprehension. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.
- Peng Cui, Dongyao Hu, and Le Hu. 2021. Listreader: Extracting list-form answers for opinion questions. *CoRR*, abs/2110.11692.
- Pradeep Dasigi, Nelson F. Liu, Ana Marasovic, Noah A. Smith, and Matt Gardner. 2019. Quoref: A reading comprehension dataset with questions requiring coreferential reasoning. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019, pages 5924–5931. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of*

the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 4171–4186. Association for Computational Linguistics.

- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 2368–2378. Association for Computational Linguistics.
- Minghao Hu, Yuxing Peng, Zhen Huang, and Dongsheng Li. 2019. A multi-type multi-span network for reading comprehension that requires discrete reasoning. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019, pages 1596– 1606. Association for Computational Linguistics.
- Yiming Ju, Weikang Wang, Yuanzhe Zhang, Suncong Zheng, Kang Liu, and Jun Zhao. 2022. CMQA: A dataset of conditional question answering with multiple-span answers. In Proceedings of the 29th International Conference on Computational Linguistics, COLING 2022, Gyeongju, Republic of Korea, October 12-17, 2022, pages 1697–1707. International Committee on Computational Linguistics.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453– 466.
- John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In Proceedings of the Eighteenth International Conference on Machine Learning (ICML 2001), Williams College, Williamstown, MA, USA, June 28 - July 1, 2001, pages 282–289. Morgan Kaufmann.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. ALBERT: A lite BERT for self-supervised learning of language representations. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.
- Kenton Lee, Tom Kwiatkowski, Ankur P. Parikh, and Dipanjan Das. 2016. Learning recurrent span representations for extractive question answering. *CoRR*, abs/1611.01436.

Seongyun Lee, Hyunjae Kim, and Jaewoo Kang. 2023. Liquid: A framework for list question answering dataset generation. In *The Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023.* 

616

617

618

625

634

636

641

642

668

672

- Haonan Li, Martin Tomko, Maria Vasardani, and Timothy Baldwin. 2022. Multispanqa: A dataset for multispan question answering. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022, pages 1250–1260. Association for Computational Linguistics.
  - Dayiheng Liu, Yeyun Gong, Jie Fu, Yu Yan, Jiusheng Chen, Daxin Jiang, Jiancheng Lv, and Nan Duan.
    2020. Rikinet: Reading wikipedia pages for natural question answering. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020, pages 6762–6771. Association for Computational Linguistics.
  - 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. *CoRR*, abs/1907.11692.
  - Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. Multi-task identification of entities, relations, and coreference for scientific knowledge graph construction. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 -November 4, 2018*, pages 3219–3232. Association for Computational Linguistics.
  - Yi Luan, Dave Wadden, Luheng He, Amy Shah, Mari Ostendorf, and Hannaneh Hajishirzi. 2019. A general framework for information extraction using dynamic span graphs. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 3036–3046. Association for Computational Linguistics.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018.
  Know what you don't know: Unanswerable questions for squad. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 2: Short Papers, pages 784–789. Association for Computational Linguistics.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100, 000+ questions for machine comprehension of text. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016, pages 2383–2392. The Association for Computational Linguistics.

Azriel Rosenfeld and Mark Thurston. 1971. Edge and curve detection for visual scene analysis. *IEEE Trans. Computers*, 20(5):562–569. 673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

709

710

711

712

713

714

716

718

719

720

721

723

724

725

726

- Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the conll-2003 shared task: Languageindependent named entity recognition. In Proceedings of the Seventh Conference on Natural Language Learning, CoNLL 2003, Held in cooperation with HLT-NAACL 2003, Edmonton, Canada, May 31 -June 1, 2003, pages 142–147. ACL.
- Elad Segal, Avia Efrat, Mor Shoham, Amir Globerson, and Jonathan Berant. 2020. A simple and effective model for answering multi-span questions. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 3074–3080. Association for Computational Linguistics.
- Min Joon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. 2017. Bidirectional attention flow for machine comprehension. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net.
- Jianlin Su, Ahmed Murtadha, Shengfeng Pan, Jing Hou, Jun Sun, Wanwei Huang, Bo Wen, and Yunfeng Liu. 2022. Global pointer: Novel efficient spanbased approach for named entity recognition. *CoRR*, abs/2208.03054.
- 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, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020,* pages 9016–9023. AAAI Press.
- Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. 2015. Pointer networks. *Advances in neural information processing systems*, 28.
- David Wadden, Ulme Wennberg, Yi Luan, and Hannaneh Hajishirzi. 2019. Entity, relation, and event extraction with contextualized span representations. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019, pages 5783–5788. Association for Computational Linguistics.
- Shuohang Wang and Jing Jiang. 2017. Machine comprehension using match-lstm and answer pointer. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net.

Junjie Yang, Zhuosheng Zhang, and Hai Zhao. 2021. Multi-span style extraction for generative reading comprehension. In Proceedings of the Workshop on Scientific Document Understanding co-located with 35th AAAI Conference on Artificial Inteligence, SDU@AAAI 2021, Virtual Event, February 9, 2021, volume 2831 of CEUR Workshop Proceedings. CEUR-WS.org.

728

729

736

737

738

740

741

742

743

744 745

746 747

748

751

753

754

755

765

770

771

772

773

774

775

776

779

780 781

783 784

- Xuchen Yao, Benjamin Van Durme, Chris Callison-Burch, and Peter Clark. 2013. Answer extraction as sequence tagging with tree edit distance. In *Human Language Technologies: Conference of the North American Chapter of the Association of Computational Linguistics, Proceedings, June 9-14, 2013, Westin Peachtree Plaza Hotel, Atlanta, Georgia, USA,* pages 858–867. The Association for Computational Linguistics.
- Deming Ye, Yankai Lin, Jiaju Du, Zhenghao Liu, Peng Li, Maosong Sun, and Zhiyuan Liu. 2020. Coreferential reasoning learning for language representation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 7170–7186. Association for Computational Linguistics.
- Wonjin Yoon, Richard Jackson, Aron Lagerberg, and Jaewoo Kang. 2022. Sequence tagging for biomedical extractive question answering. *Bioinform.*, 38(15):3794–3801.
- Adams Wei Yu, David Dohan, Minh-Thang Luong, Rui Zhao, Kai Chen, Mohammad Norouzi, and Quoc V. Le. 2018. Qanet: Combining local convolution with global self-attention for reading comprehension. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net.
- Juntao Yu, Bernd Bohnet, and Massimo Poesio. 2020. Named entity recognition as dependency parsing. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020, pages 6470–6476. Association for Computational Linguistics.
- Zheng Yuan, Chuanqi 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, Dublin, Ireland, May* 22-27, 2022, pages 3174–3186. Association for Computational Linguistics.
- Zhuosheng Zhang and Hai Zhao. 2022. Tracing origins: Coreference-aware machine reading comprehension. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022, pages 1281–1292. Association for Computational Linguistics.
- Ming Zhu, Aman Ahuja, Da-Cheng Juan, Wei Wei, and Chandan K. Reddy. 2020. Question answering

with long multiple-span answers. In Findings of the<br/>Association for Computational Linguistics: EMNLP7862020, Online Event, 16-20 November 2020, volume<br/>EMNLP 2020 of Findings of ACL, pages 3840–3849.789Association for Computational Linguistics.790