Context-Aware Transformer Pre-Training for Answer Sentence Selection

Anonymous ARR submission

Abstract

Answer Sentence Selection (AS2) is a core component for building an accurate Question Answering pipeline. AS2 models rank a set of 004 candidate sentences based on how likely they answer a given question. The state of the art in AS2 exploits pre-trained transformers by trans-007 ferring them on large annotated datasets, while using local contextual information around the 009 candidate sentence. In this paper, we propose three pre-training objectives designed to mimic the downstream fine-tuning task of contextual AS2. This allows for specializing LMs when 013 fine-tuning for contextual AS2. Our experiments on three public and two large-scale in-015 dustrial datasets show that our pre-training approaches (applied to RoBERTa and ELECTRA) 017 can improve baseline contextual AS2 accuracy by up to 8% on some datasets.

1 Introduction

021

034

040

Answer Sentence Selection (AS2) is a fundamental task in QA, which consists of re-ranking a set of answer sentence candidates according to how correctly they answer a given question. From a practical standpoint, AS2-based QA systems can operate under much lower latency constraints than corresponding Machine Reading (MR) based QA systems. This is because AS2 systems process several sentences/documents in parallel, while MR systems parse the entire document/passage in a sliding window fashion before finding the answer span (Garg and Moschitti, 2021).

Modern AS2 systems (Garg et al., 2020; Laskar et al., 2020) use transformers to cross-encode question and answer candidates together. Recently, Lauriola and Moschitti (2021) proved that performing answer ranking using only the candidate sentence is sub-optimal, for e.g., the answer sentence may contain unresolved coreference with entities, or the sentence may lack specific context for answering the question. Several works (Ghosh et al., 2016; Tan et al., 2018; Han et al., 2021) have explored performing AS2 using context around answer candidates (for example, adjacent sentences) towards improving performance. Local contextual information, i.e., the previous and next sentences of the answer candidates, can help coreference disambiguation, and provide additional knowledge to the model. This helps to rank the best answer at the top, with minimal increase in compute requirements. 042

043

044

045

046

047

051

054

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

081

Previous research works (Lauriola and Moschitti, 2021; Han et al., 2021) have directly used existing pre-trained transformer encoders for contextual AS2, by fine-tuning them on an input comprising of multiple sentences with different roles, i.e., the question, answer candidate, and context (previous and following sentences around the candidate). This structured input creates practical challenges during fine-tuning, as standard pre-training approaches do not align well with the downstream contextual AS2 task, e.g., the language model does not know the role of each of these multiple sentences in the input. In other words, the extended sentence-level embeddings have to be learnt directly during fine-tuning, causing underperformance empirically. This effect is amplified when the downstream data for fine-tuning is small, indicating models struggling to exploit the context.

In this paper, we tackle the aforementioned issues by designing three pre-training objectives that structurally align with the final contextual AS2 task, and can help improve the performance of language models when fine-tuned for AS2. Our pre-training objectives exploit information in the structure of paragraphs and documents to pre-train the context slots in the transformer text input. We evaluate our strategies on two popular pre-trained transformers over five datasets. The results show that our approaches using structural pre-training can effectively adapt transformers to process contextualized input, improving accuracy by up to 8% when compared to the baselines on some datasets. We plan to release our code and pre-trained models.

2 Related Work

084

086

095

096

098

100

101

102

129

130

131

132

133

Answer Sentence Selection: TANDA (Garg et al., 2020) established the SOTA for AS2 using a large dataset (ASNQ) for transfer learning. Other approaches for AS2 include: separate encoders for question and answers (Bonadiman and Moschitti, 2020), and compare-aggregate and clustering to improve answer relevance ranking (Yoon et al., 2019). Contextual AS2: Ghosh et al. (2016) use LSTMs for answers and topics, improving accuracy for next sentence selection. Tan et al. (2018) use GRUs to model answers and local context, improving performance on two AS2 datasets. Lauriola and Moschitti (2021) propose a transformer encoder that uses context to better disambiguate between answer candidates. Han et al. (2021) use unsupervised similarity matching techniques to extract relevant context for answer candidates from documents. Refer to Appendix **B** for a discussion on different forms of contextual information for AS2.

Pre-training Objectives: Pre-training sentence-103 level objectives such as NSP (Devlin et al., 2019) 104 and SOP (Lan et al., 2020) have been widely ex-105 plored for transformers to improve accuracy for 106 downstream classification tasks. However, the ma-107 jority of these objectives are agnostic of the final 108 tasks. End task-aware pre-training has been studied 109 for summarization (Rothe et al., 2021), dialogue 110 (Li et al., 2020), passage retrieval (Gao and Callan, 111 2021), MR (Ram et al., 2021) and multi-task learn-112 ing (Dery et al., 2021). Lee et al. (2019), Chang 113 et al. (2020a) and Sachan et al. (2021) use the In-114 verse Cloze task to improve retrieval performance 115 for bi-encoders, by exploiting paragraph structure 116 via self-supervised objectives. For AS2, recently 117 Di Liello et al. (2022a) proposed paragraph-aware 118 pre-training for joint classification of multiple can-119 didates. Di Liello et al. (2022b) propose a sentence-120 level pre-training paradigm for AS2 by exploiting 121 document and paragraph structure. However, these 122 works do not consider the structure of the downstream task (specifically contextual AS2). To the 124 best of our knowledge, ours is the first work to 125 study transformer pre-training strategies for AS2 126 augmented with context using cross-encoders. 127

3 Contextual AS2

AS2: Given a question q and a set of answer candidates $S = \{s_1, \ldots, s_n\}$, the goal is to find the best s_k that answers q. This is typically done by learning a binary classifier C of answer correctness by independently feeding the pairs $(q, s_i), i \in$ $\{1, ..., n\}$ as input to **C**, and making **C** predict whether s_i correctly answers q or not. At inference time, we find the best answer for q by selecting the answer candidate s_k which scores the highest probability of correctness $k = \arg \max_i \mathbf{C}(q, s_i)$. **Contextual AS2**: Contextual models for AS2 exploit additional context to improve the final accuracy. This has been shown to be effective (Lauriola and Moschitti, 2021) in terms of overcoming coreference disambiguation and lack of enough information to rank the best answer at the top. Different from the above case, contextual AS2 models receive as input a tuple (q, s_i, c_i) where c_i is the additional context. c_i is usually the sentences immediately before and after the answer candidate.

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

4 Context-aware Pre-training Objectives

We design a transformer pre-training task that aligns well with fine-tuning contextual AS2 models, both *structurally* and *semantically*. We exploit the division of large corpora in documents and the subdivision of documents in paragraphs as a source of supervision. We provide triplets of text spans (a, b, c) as model inputs when pre-training, which emulates the structure of (q, s_i, c_i) for contextual AS2 models, where a, b and c play the analogous role of the question, the candidate sentence (that needs to be classified), and the context (which helps in predicting (a, b) correctness), respectively. Formally, given a document D from the pre-training corpus, the task is to infer if a and b are two sentences extracted from the same paragraph $P \in D$.

The semantics learned by connecting sentences in the same paragraph transfer well downstream, as the model can re-use previously learned relations between entities and concepts, and apply them between question and answer candidates. Relations in one sentence may be used to formulate questions that can be answered in the other sentence, which is most likely to happen for sentences in the same paragraph. We expand this discussion with some examples in Appendix A. We term this task: "Sentences in Same Paragraph (SSP)" and design three ways of choosing the appropriate contextual information c. We present details on how we sample spans a, b and c from the pre-training documents. Static Document-level Context (SDC) Here, we choose the context c to be the first paragraph P_0 of $D = \{P_0, ..., P_n\}$ from which b is extracted. This is based on the intuition that the first paragraph acts as a summary of a document's content (Chang et al., 2020a): this strong context can help

the model at identifying if b is extracted from the 185 same paragraph as a. We call this static document-186 level context since the contextual information c is constant for any b extracted from the same document D. Specifically, the positive examples are created by sampling a and b from a single ran-190 dom paragraph $P_i \in D, i > 0$. For the previously 191 chosen a, we create hard negatives by randomly 192 sampling a sentence b from different paragraphs 193 $P_j \in D, j \neq i \land j > 0$. We set $c = P_0$ for this 194 negative example as well since b still belongs to 195 D. We create easy negatives for a chosen a by 196 sampling b from a random paragraph P'_i in another 197 document $D' \neq D$. In this case, c is chosen as the 198 first paragraph P'_0 of D' since the context in the 199 downstream AS2 task is associated with the answer candidate, and not with the question.

202

206

210

211

212

213

214

Dynamic Paragraph-level Context (DPC) We dynamically select the context c to be the paragraph from which the sentence b is extracted. We create positive examples by sampling a and b from a single random paragraph $P_i \in D$, and we set the context as the remaining sentences in P_i , i.e., $c = P_i \setminus \{a, b\}$. Note that leaving a and b in P_i would make the task trivial. For the previously chosen a, we create hard negatives by sampling b from another random paragraph $P_j \in D, j \neq i$, and setting $c = P_j \setminus \{b\}$. We create easy negatives for a chosen a by sampling b from a random P'_i in another document $D' \neq D$, and setting $c = P'_i \setminus \{b\}$.

Dynamic Sentence-level Local Context (DSLC) 215 216 We choose c to be the local context around the sentence b, i.e, the concatenation of the previous and 217 next sentence around b in $P \in D$. To deal with 218 corner cases, we require at least one of the previ-219 ous or next sentences of b to exist (e.g., the next sentence may not exist if b is the last sentence of the paragraph P). We term this DSLC as the con-222 textual information c is specified at sentence-level and changes correspondingly to every sentence b extracted from D. We create positive pairs similar to SDC and DPC by sampling a and b from the 226 same paragraph $P_i \in D$, with c being the local 227 context around b in P_i (and $a \notin c$). We automatically discard paragraphs that are not long enough to ensure the creation of a positive example. We generate hard negatives by sampling b from another $P_i \in D, j \neq i$, while for easy negatives, we sample b from a $P'_i \in D', D' \neq D$ (in both cases c is set as the local context around b).

5 Datasets

Pre-Training To perform a fair comparison and avoid any improvement stemming from additional pre-training data, we use the same corpora as RoBERTa (Liu et al., 2019). This includes the English Wikipedia, the BookCorpus (Zhu et al., 2015), OpenWebText (Gokaslan and Cohen, 2019) and CC-News. See Appendix C.1 for more details. We transform the datasets to implement the pretraining objectives that we described in Section 4. **Contextual AS2** We evaluate our pre-trained models on three public and two industrial datasets for contextual AS2. For all datasets, we use the standard "clean" setting, by having at least one positive and one negative candidate per question in the dev. and test sets. We measure performance using Precision-at-1 (P@1) and Mean Average Precision (MAP). Datasets statistics and details are presented in Appendix C.2.

• ASNQ is a large scale AS2 dataset (Garg et al., 2020) derived from NQ (Kwiatkowski et al., 2019). The questions are user queries from Google search, and answers are extracted from Wikipedia.

• WikiQA is a small dataset (Yang et al., 2015) for AS2 with questions extracted from Bing search engine and answer candidates retrieved from the first paragraph of Wikipedia articles.

• NewsAS2 is a large AS2 dataset created from NewsQA (Trischler et al., 2017), a MR dataset, following the procedure of Garg et al. for ASNQ. The dataset contains ~70K human generated questions with answers extracted from *CNN/Daily Mail*.

• IQAD is a large scale industrial dataset containing de-identified questions asked by users to a popular commercial virtual assistant. IQAD contains \sim 220k questions where answers are retrieved from a large web index (~1B web pages) using Elasticsearch. We use two different evaluation benchmarks for IQAD: (i) IQAD Bench 1, which contains 2.2k questions with \sim 15 answer candidates annotated for correctness by crowd workers and (ii) IQAD Bench 2, which contains 2k questions with \sim 15 answer candidates annotated with explicit fact verification guidelines for correctness by crowd workers. (Our manual analysis indicates a higher annotation quality for QA pairs in Bench 2 than Bench 1). Results on IQAD are presented relative to a baseline due to the data being internal.

6 Experiments

Continuous Pre-Training We use RoBERTa-Base and ELECTRA-Base public checkpoints (pre-

Model	Context	ASNQ		WikiQA		NewsAS2		IQAD Bench 1		IQAD Bench 2	
		MAP	P@1	MAP	P@1	MAP	P@1	MAP	P@1	MAP	P@1
ELECTRA-Base	x	69.3 (0.0)	65.0 (0.2)	85.7 (0.9)	78.5 (1.6)	81.3 (0.2)	75.6 (0.2)	Base	eline	Base	eline
ELECTRA-Base 🜲	1	72.3 (0.6)	68.1 (0.8)	83.1 (1.3)	73.8 (2.1)	82.0 (0.4)	76.0 (0.5)	-0.6%	-1.0%	-0.4%	-0.9%
(Ours) ELECTRA-Base + SSP (SDC)	1	<u>74.7</u> (0.5)	<u>69.6</u> (0.3)	88.7 (0.1)	82.9 (0.2)	<u>82.7</u> (0.2)	77.0 (0.4)	+1.2%	+0.6%	+0.9%	+1.4%
(Ours) ELECTRA-Base + SSP (DPC)	1	<u>74.4</u> (0.2)	70.5 (0.2)	<u>88.0</u> (0.6)	81.3 (0.6)	82.7 (0.5)	<u>77.3</u> (0.7)	+0.4%	-0.6%	+0.4%	+0.1%
(Ours) ELECTRA-Base + SSP (DSLC)	1	<u>74.3</u> (0.3)	<u>70.0</u> (0.8)	<u>87.0</u> (0.9)	<u>79.7</u> (1.4)	82.8 (0.4)	77.3 (0.5)	+1.0%	+0.6%	+0.2%	0.0%
(Ours) ELECTRA-Base + SSP (All)	1	<u>73.8</u> (0.4)	68.8 (0.4)	<u>87.5</u> (0.5)	<u>81.5</u> (0.7)	<u>82.7</u> (0.2)	<u>77.2</u> (0.3)	+0.1%	-0.4%	+0.1%	-0.1%
RoBERTa-Base	x	68.2 (0.5)	63.5 (0.5)	85.1 (1.9)	77.2 (3.1)	81.7 (0.1)	76.2 (0.2)	+0.6%	+0.1%	+0.7%	+1.3%
RoBERTa-Base 🌲	1	71.6 (0.6)	67.6 (0.6)	84.4 (1.5)	77.0 (2.1)	82.4 (0.2)	76.6 (0.7)	+0.4%	0.0%	+1.1%	+1.7%
(Ours) RoBERTa-Base + SSP (SDC)	1	<u>73.1</u> (0.5)	68.7 (0.8)	<u>87.8</u> (0.6)	<u>81.8</u> (0.9)	<u>82.8</u> (0.1)	76.9 (0.2)	<u>+1.7%</u>	+3.0%	+1.0%	+1.7%
(Ours) RoBERTa-Base + SSP (DPC)	1	73.2 (0.4)	69.2 (0.5)	89.9 (0.2)	85.2 (0.4)	82.3 (0.1)	76.0 (0.1)	+0.4%	+1.2%	+1.2%	+2.7%
(Ours) RoBERTa-Base + SSP (DSLC)	1	<u>72.9</u> (0.4)	<u>69.0</u> (0.3)	<u>87.8</u> (0.9)	81.6 (1.3)	82.6 (0.2)	77.0 (0.2)	+0.6%	+1.5%	+1.0%	+1.4%
(Ours) RoBERTa-Base + SSP (All)	1	72.9 (0.6)	68.2 (0.8)	<u>88.2</u> (0.9)	<u>82.4</u> (1.7)	<u>83.0</u> (0.2)	77.3 (0.5)	<u>+1.2%</u>	<u>+2.4%</u>	+1.4%	+2.2%

Table 1: Results (std. dev. in parenthesis) on AS2. Models with \clubsuit are from (Lauriola and Moschitti, 2021). \checkmark and \checkmark denote whether local contextual information was used in fine-tuning. SDC, DPC and DSLC indicate the pre-training variants of the SSP task that we propose. Best results are in bold while we underline statistically significant improvements over the two contextual baselines (\clubsuit) using a Student *t*-test with 95% of confidence level.

training from scratch would have required large amounts of computational resources), and perform continuous pre-training using our objectives for $\sim 10\%$ of the compute used by the original models. Complete details are given in Appendix E. We experiment with each of our pre-training objectives independently, as well as combining all of them.

287

290

296

297

319

320

Fine-Tuning We fine-tune each continuously pretrained model on all the AS2 datasets. As baselines, we consider (i) standard pairwise-finetuned AS2 models, using only the question and the answer candidate, and (ii) contextual fine-tuned AS2 models from (Lauriola and Moschitti, 2021), which use the question, answer candidate and local context.

Results Table 1 summarizes the results of our 301 experiments averaged across 5 runs. On ASNQ, 302 our pre-trained models get 3.8 - 5.5% improvement in P@1 over the baseline using only the question and answer. Our models also outper-304 form the stronger contextual AS2 baselines (1.6% 305 306 with RoBERTa and 2.4% with ELECTRA), indicating that our task-aware pre-training can help im-307 prove the downstream fine-tuning performance. On NewsAS2, we observe a similar trend, where all our models (except one) outperform both the standard 310 and contextual baselines. On WikiQA, a smaller 311 dataset, the contextual baseline underperforms the 312 non-contextual baseline, highlighting that with few 313 samples the model struggles to adapt and reason over three text spans. Our pre-training approaches 315 provide the maximum performance improvement 316 on WikiQA (up to 8 - 9.1% improvement over the 317 non-contextual and contextual baselines). 318

On IQAD, we observe that the contextual baseline performs on par or lower than the noncontextual baseline, indicating that off-the-shelf transformers cannot effectively exploit the context available for this dataset. The answer candidates and context for IQAD are extracted from millions of web documents. Thus, learning from the context in IQAD is a harder task than learning from it on ASNQ, where the context belongs to a single Wikipedia document. Our pre-trained models help to process the diverse and possibly noisy context of IQAD, and produce a significant improvement in P@1 over the contextual baseline.

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

347

348

350

351

352

353

354

355

356

357

358

The DPC and DSLC approaches align well (often having overlapping or identical contexts for the same (a, b) input): this explains their comparable performance across all datasets. In SDC, the context c can potentially be very different from (a, b), and this may help in exploiting information from multiple documents/domains as in case of IQAD. For these reasons, we believe DPC and DSLC should be used when answer candidates are extracted from the same document, while SDC works best with candidates collected across multiple documents. We present an extended discussion of our results in Appendix G. Also, we observe that combining all the objectives together does not always outperform the individual objectives, which is probably due to the misalignment between the different approaches for sampling context in our pre-training strategies.

7 Conclusions

In this paper, we have proposed three pre-training strategies for transformers, which (i) are aware of the downstream task of contextual AS2, and (ii) use the document and paragraph structure information to define effective objectives. Our experiments on three public and two industrial datasets using two transformer models show that our pre-training strategies can provide significant improvement over the contextual AS2 models.

References

366

367

374

376

377

380

384

387

400

401

402

403

404

405

406

407

408

409

410

411

412

413

- Daniele Bonadiman and Alessandro Moschitti. 2020. A study on efficiency, accuracy and document structure for answer sentence selection. In Proceedings of the 28th International Conference on Computational Linguistics, pages 5211–5222, Barcelona, Spain (Online). International Committee on Computational Linguistics.
 - Wei-Cheng Chang, Felix X. Yu, Yin-Wen Chang, Yiming Yang, and Sanjiv Kumar. 2020a. Pre-training tasks for embedding-based large-scale retrieval.
 - Wei-Cheng Chang, Felix X. Yu, Yin-Wen Chang, Yiming Yang, and Sanjiv Kumar. 2020b. Pre-training tasks for embedding-based large-scale retrieval. In International Conference on Learning Representations.
 - Lucio M. Dery, Paul Michel, Ameet Talwalkar, and Graham Neubig. 2021. Should we be pre-training? an argument for end-task aware training as an alternative.
 - Nicki Skafte Detlefsen, Jiri Borovec, Justus Schock, Ananya Harsh Jha, Teddy Koker, Luca Di Liello, Daniel Stancl, Changsheng Quan, Maxim Grechkin, and William Falcon. 2022. Torchmetrics - measuring reproducibility in pytorch. *Journal of Open Source Software*, 7(70):4101.
 - 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.
 - Luca Di Liello, Siddhant Garg, Luca Soldaini, and Alessandro Moschitti. 2022a. Paragraph-based transformer pre-training for multi-sentence inference. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics, Seattle, Washington. Association for Computational Linguistics.
 - Luca Di Liello, Siddhant Garg, Luca Soldaini, and Alessandro Moschitti. 2022b. Pre-training transformer models with sentence-level objectives for answer sentence selection.
 - William Falcon et al. 2019. Pytorch lightning. GitHub. Note: https://github.com/PyTorchLightning/pytorchlightning, 3(6).
 - Luyu Gao and Jamie Callan. 2021. Unsupervised corpus aware language model pre-training for dense passage retrieval.
- Siddhant Garg and Alessandro Moschitti. 2021. Will this question be answered? question filtering via

answer model distillation for efficient question answering. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7329–7346, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. 414

415

416

417

418

419

420

421

422

423

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

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

- Siddhant Garg, Thuy Vu, and Alessandro Moschitti. 2020. Tanda: Transfer and adapt pre-trained transformer models for answer sentence selection. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):7780–7788.
- Shalini Ghosh, Oriol Vinyals, Brian Strope, Scott Roy, Tom Dean, and Larry Heck. 2016. Contextual lstm (clstm) models for large scale nlp tasks.
- Aaron Gokaslan and Vanya Cohen. 2019. Openwebtext corpus. http://Skylion007.github.io/ OpenWebTextCorpus.
- Rujun Han, Luca Soldaini, and Alessandro Moschitti. 2021. Modeling context in answer sentence selection systems on a latency budget.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics*.
- 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.
- Md Tahmid Rahman Laskar, Jimmy Xiangji Huang, and Enamul Hoque. 2020. Contextualized embeddings based transformer encoder for sentence similarity modeling in answer selection task. In *Proceedings of the 12th Language Resources and Evaluation Confer ence*, pages 5505–5514, Marseille, France. European Language Resources Association.
- Ivano Lauriola and Alessandro Moschitti. 2021. Answer sentence selection using local and global context in transformer models. In *ECIR 2021*.
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. Latent retrieval for weakly supervised open domain question answering.
- Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario Šaško, Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Clément Delangue, Théo Matussière, Lysandre Debut, Stas Bekman, Pierric Cistac, Thibault Goehringer, Victor Mustar, François

548

549

526

Lagunas, Alexander Rush, and Thomas Wolf. 2021. Datasets: A community library for natural language processing. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 175–184, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

470

471

472

473

474

475

476

477

478

479

480

481

482

483 484

485

486

487

488

489

490

491

492

493

494 495

496

497

498 499

505 506

507

508

509

510 511

512

513

514

515

516

517

518

519 520

521

524

- Junlong Li, Zhuosheng Zhang, Hai Zhao, Xi Zhou, and Xiang Zhou. 2020. Task-specific objectives of pretrained language models for dialogue adaptation.
- 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.
 - Edward Loper and Steven Bird. 2002. Nltk: The natural language toolkit. *CoRR*, cs.CL/0205028.
 - Ori Ram, Yuval Kirstain, Jonathan Berant, Amir Globerson, and Omer Levy. 2021. Few-shot question answering by pretraining span selection. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 3066–3079, Online. Association for Computational Linguistics.
 - Sascha Rothe, Joshua Maynez, and Shashi Narayan. 2021. A thorough evaluation of task-specific pretraining for summarization. In *Proceedings of the* 2021 Conference on Empirical Methods in Natural Language Processing, pages 140–145, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
 - Devendra Singh Sachan, Mostofa Patwary, Mohammad Shoeybi, Neel Kant, Wei Ping, William L Hamilton, and Bryan Catanzaro. 2021. End-to-end training of neural retrievers for open-domain question answering.
 - Chuanqi Tan, Furu Wei, Qingyu Zhou, Nan Yang, Bowen Du, Weifeng Lv, and Ming Zhou. 2018. Context-aware answer sentence selection with hierarchical gated recurrent neural networks. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 26(3):540–549.
 - Trieu H. Trinh and Quoc V. Le. 2018. A simple method for commonsense reasoning. *ArXiv*, abs/1806.02847.
- Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordoni, Philip Bachman, and Kaheer Suleman. 2017. NewsQA: A machine comprehension dataset. In Proceedings of the 2nd Workshop on Representation Learning for NLP, pages 191–200, Vancouver, Canada. Association for Computational Linguistics.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen,

Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

- Yi Yang, Scott Wen-tau Yih, and Chris Meek. 2015. Wikiqa: A challenge dataset for open-domain question answering. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. ACL - Association for Computational Linguistics.
- Seunghyun Yoon, Franck Dernoncourt, Doo Soon Kim, Trung Bui, and Kyomin Jung. 2019. A compareaggregate model with latent clustering for answer selection.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *The IEEE International Conference on Computer Vision (ICCV)*.

Appendix

550

551

552

554

555

556

557

564

565

566

569

570

571

575

577

578

582

588

590

594

595

596

A Intuition behind SSP Pre-training

In this Section, we expand on the discussion of why the task of predicting whether 2 sentences are extracted from same paragraph of a document (SSP) can help the downstream task of AS2. We start by considering an example of a paragraph extracted from Wikipedia composed of three sentences:

 s_1 : Lovato was brought up in Dallas, Texas; she began playing the piano at age seven and guitar at ten, when she began dancing and acting classes.

 s_2 : In 2002, Lovato began her acting career on the children's television series Barney & Friends, portraying the role of Angela.

 s_3 : She appeared on Prison Break in 2006 and on Just Jordan the following year.

Given a question of the type "What are the acting roles of X", a standard language model can easily reason to select answers of the type "Xacted/played in Y", as this is about matching the subject argument of the question with the object argument of the answer, for the same predicate acting/playing. However, the same LM would have a harder time selecting answers of the type "X appeared in Y" because this requires learning the relation between the entire predicate argument structure of acting vs. the one of appearing.

In contrast, a LM pre-trained using our SSP approach can learn these implications, as it reasons about these concepts and relations from s_3 , e.g., "appearing in Prison Break and Just Jordan" (which are TV series), are related to concepts and relations from s_2 , e.g., "having an acting career". This is learned as our model reasons to connect sentences which are in the same paragraphs. During fine-tuning, we use the question text in place of one of the two sentences but the model can still reuse the relations learned between sentences, and apply them between question and the answer candidates.

Speaking in more general terms, some of the relations that are in one sentence may be used to formulate questions that can be answered in the other sentences. This happens with the high probability for sentences appearing within the same paragraph. This rationale also motivates our choice of easy and hard negative examples.

B Contextual Information for AS2

In addition to local context around answer candidates (the previous and successive sentences), other contextual signals can also be incorporated to improve the relevance ranking of answer candidates. Meta-information like document title, abstract/firstparagraph, domain name, etc. corresponding to the document containing the answer candidates can help answer ranking. These signals differ from the previously mentioned local answer context as they provide "global" contextual information pertaining to the documents for AS2. 599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

Lauriola and Moschitti (2021) present an initial exploration of global contextual signals for AS2, along with their experiments on local answer context. This is done by computing a Bag of Words (BoW) vector representation over the document from which the answer is extracted, and then concatenating this with the [CLS] embeddings from the transformer before feeding it to the final classification layer (the global context is not encoded as a transformer input). Empirically, they show that using local answer context outperforms global answer context on some datasets, while achieves comparable performance on others.

In this work, our primary objective is to design pre-training techniques that help the transformer encoder reason over an additional input (the context) which is related to the answer candidate in some manner to eventually help improve the relative ranking of answer candidates. We don't put any constraints on what form this context can take, for eg: it can be derived from the same paragraph as the answer candidate, or it can be derived from the first paragraph of the document, etc. Due to global contextual information not being available for all the datasets we consider (documents having missing paragraphs in some cases), we empirically evaluate our approaches using only the local answer context. Our Static Document-level Context (SDC) objective, which uses the first paragraph of the document for the context input slot, captures global information pertaining to the document ((Chang et al., 2020b) show that the first paragraph acts as a summary of a document's content). We hypothesize that this will improve downstream performance using other global contextual signals in addition to local answer context.

C Datasets

C.1 Pre-training

RoBERTa was trained over English Wikipedia, the BookCorpus (Zhu et al., 2015), STORIES (Trinh and Le, 2018), OpenWebText (Gokaslan and Co-

- 674
- 675

679

681

hen, 2019) and CC-News. However, STORIES is no longer publicly available, and thus we ignore it.

We preprocess Wikipedia, the BookCorpus, CC-News and OpenWebText by filtering away: (i) sentences having a length smaller than 20 characters, (ii) paragraphs shorter than 60 characters and (iii) documents shorter than 200 characters. We split paragraphs in sequences of sentences using the NLTK tokenizer (Loper and Bird, 2002) and we create the datasets for continuous pre-training following the definitions in Section 4.

For each objective, we sample randomly up to 2 hard negatives and additional easier negatives until the total number is 4. Instead of reasoning in terms of sentences, we designed our objectives to create a and b as small spans composed of 1 or more contiguous sentences. For a, we keep the length equal to 1 sentence because it emulates the question, which usually is just a single sentence. For b, we randomly sample the length between 1 and 3. The length of the context c cannot be decided a priori because it depends on the specific pre-training objective and the length of the paragraph.

All the resulting continuous pre-training datasets are about 300GB in size (uncompressed) and contain around 350M training examples each.

C.2 Fine-Tuning

The statistics on the number of unique questions and question-answer pairs for each fine-tuning dataset are provided in Table 2. While ASNQ has a huge number of negatives for each question (more than 300 on average), the other datasets have a smaller number of answer candidates per question. All datasets are converted to the "clean" setting, which means questions without at least a positive and a negative answer candidates are removed, which is a standard practice in AS2.

NewsAS2 was created by splitting each document in NewsQA into individual sentences with 687 the NLTK tokenizer (Loper and Bird, 2002). Then, for each sentence, we assigned a positive label if it contained at least one of the annotated answers for that document, a negative label otherwise. This lead to a dataset with 1.69% positives sentences per query in the training set, 1.66% in the dev set and 693 1.68% in the test set. We will release this NewsAS2 dataset along with code and models from our paper.

Dataset	Т	rain]	Dev	Test		
	#Q	#QA	#Q	#QA	#Q	#QA	
ASNQ	57242	20377568	1336	463914	1336	466148	
WikiQA	2118	20360	122	1126	237	2341	
IQAD	221334	3894129	2434	43369	2252 2088	38587 33498	
NewsAS2	71561	1840533	2102	51844	2083	51472	

Table 2: Number or unique questions and questionanswer pairs in the fine-tuning datasets. IQAD Bench 1 and Bench 2 sizes are mentioned in the Test set column corresponding to IQAD.

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

Frameworks & Infrastructure D

Our framework is based on (i) HuggingFace Transformers (Wolf et al., 2020) for model architecture, (ii) HuggingFace Datasets (Lhoest et al., 2021) for data processing, (iii) PyTorch-Lightning for distributed training (Falcon et al., 2019) and (iv) TorchMetrics for AS2 evaluation metrics (Detlefsen et al., 2022).

We performed our pre-training experiments for every model on 8 NVIDIA A100 GPUs with 40GB of memory each, using fp16 for tensor core acceleration.

Е **Details of Continuous Pre-Training**

We experiment with RoBERTa-Base and ELECTRA-Base public checkpoints. RoBERTa-Base contains 124M parameters while ELECTRA-Base contains 33M parameters in the generator and 108M in the discriminator.

We do continuous pre-training starting from the aforementioned models for 400K steps with a batch size of 4096 examples and a triangular learning rate with a peak value of 10^{-4} and 10K steps of warm-up. In order to save resources, we found it beneficial to reduce the maximum sequence length to 128 tokens. In this setting, our models see about 210B additional tokens each, which are exactly the 10% of those used in the original RoBERTa pre-training. Moreover, in terms of complexity our objectives are more efficient because the attention computational complexity grows quadratically in the sequence length, which in our case is 4 times smaller.

We use cross-entropy as the loss function for all our pre-training and fine-tuning experiments. Specifically, for RoBERTa pre-training we sum the MLM and our proposed binary classification losses with equal weights (1.0). For ELECTRA pre-training, we sum three losses: MLM loss with a weight of 1.0, the Token Detection loss with a weight of 50.0, and our proposed binary classifica-

Model	Hyper-parameter	ASNQ	WikiQA	NewsAS2	IQAD
RoBERTa	Batch size	2048	32	256	256
	Peak LR	1e-05	5e-06	5e-06	1e-05
	Warmup steps	10K	1K	5K	5K
	Epochs	6	30	8	10
ELECTRA	Batch size	1024	128	128	256
	Peak LR	1e-05	2e-05	1e-05	2e-05
	Warmup steps	10K	1K	5K	5K
	Epochs	6	30	8	10

Table 3: Hyper-parameters used to fine-tune RoBERTa and ELECTRA on the AS2 datasets. The best hyperparameters has been chosen based on the MAP results on the validation set.

tion losses with a weight of 1.0.

During continuous pre-training, we feed the text tuples (a, b, c) (as described in Section 4) as input to the model in the following format: '[CLS]*a*[SEP]*b*[SEP]*c*[SEP]'. To provide independent sentence/segment ids to each of the inputs *a*, *b* and *c*, we initialize the sentence embeddings layers of RoBERTa and ELECTRA from scratch, and extend them to an input size of 3.

The pre-training of every model obtained by combining ELECTRA and RoBERTa architectures with our contextual pre-training objectives took around 3.5 days each on the machine configuration described in Appendix D. All the dataset preparation required 10 hours over 64 CPU cores.

F Details of Fine-Tuning

The most common paradigm for AS2 fine-tuning is to consider publicly available pre-trained transformer checkpoints (pre-trained on large amounts of raw data) and fine-tune them on the AS2 datasets. Using our proposed pre-training objectives, we are proposing stronger model checkpoints ¹ which can improve over the standard public checkpoints, and can be used as the initialization for downstream fine-tuning for contextual AS2.

To fine-tune our models on the downstream AS2 datasets, we found it is beneficial to use a very large batch size for ASNQ and a smaller one for IQAD, NewsAS2 and WikiQA. Moreover, for every experiment we used a triangular learning rate scheduler and we did early stopping on the development set if the MAP did not improve for 5 times in a row. We fixed the maximum sequence length to 256 tokens in every run, and we repeated them 3 times with different initial random seeds. We did not use weight decay but we clipped gradients larger than 1.0 in absolute value. More specifically, for the learning rate we tried all values in $\{5 * 10^{-6}, 10^{-5}, 2 * 10^{-5}\}$ for RoBERTa and in $\{10^{-5}, 2 * 10^{-5}, 5 * 10^{-5}\}$ for ELECTRA. Regarding the batch size, we tried all values in $\{512, 1024, 2048, 4096\}$ for ASNQ, in $\{64, 128, 256, 512\}$ for IQAD and NewsAS2 and in $\{16, 32, 64, 128\}$ for WikiQA. More details about final setting are given in Table 3.

774

775

776

777

778

779

780

781

782

783

784

785

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

For the pair-wise models, we format inputs as '[CLS]q[SEP] s_i [SEP]', while for contextual models we build inputs of the form '[CLS]q[SEP] s_i [SEP] c_i [SEP]'.

We do not use extended sentence/segment ids for the non-contextual baselines and retain the original model design: (i) disabled segment ids for RoBERTa and (ii) only using 2 different sentence/segment ids for ELECTRA. For the finetuning of our continuously pre-trained models as well as the contextual baseline, we use three different sentence ids corresponding to q, s and c for both RoBERTa and ELECTRA.

Finally, differently from pre-training, in finetuning we always provide the previous and the next sentence as context for a given candidate.

The contextual fine-tuning of every models on ASNQ required 6 hours per run on the machine configuration described in Appendix D. For the other fine-tuning datasets, we used a single GPU for every experiment, and runs took less than 2 hours.

G Additional Discussion of Results

In this Section, we explain the difference in performance we observe from our three pre-training objectives on different AS2 datasets. The AS2 datasets we consider for our experiments have significantly different structures: specifically, ASNQ and NewsAS2 have answer candidates being extracted from a single document (Wikipedia and CNN Daily Mail article respectively), while IQAD has answer candidates being extracted from multiple documents. This also results in the context for the former being more homogeneous (context for all candidates for a question is extracted from the same document), while for the latter the context is more heterogeneous (extracted from multiple documents for different answer candidates).

Our DPC and DSLC pre-training approaches are well aligned in terms of the context that is used to help the SSP predictions. The former uses the remainder of the paragraph P as context (after removing a and b), while the latter uses the sentence previous and next to b in P. We observe empiri-

771

772

773

736

¹We plan to release our code and pre-trained model checkpoints after the anonymity period.

cally that the contexts for DPC and DSLC often overlap partially, and are sometimes even identical (considering average length of paragraphs in the pre-training corpora is 4 sentences). This explains why models pre-trained using both these approaches perform comparably in Table 1 (with only a very small gap in P@1 performance).

> On IQAD, we observe that the SDC approach of providing context for SSP outperforms the DPC and DSLC approaches for pre-training. In SDC, the context c can potentially be very different from a and b (as it corresponds to the first paragraph of the document), and this can aid exploiting information and effectively ranking answer candidates from multiple documents (possibly from different domains) like for IQAD.

H Qualitative Examples

In Table 4 we show a comparison of the ranking produced by our models and that by the contextual baselines on some questions selected from the ASNQ test set.

ELECTRA

824

825

830

833

834

835

836

837

842

843

- Q how many games does a team have to win for the world series
 A₁ Seven games were played, with the Astros victorious after game seven, played in Los Angeles.
- A2 In 1985, the format changed to best-of-seven.
- A₃ Since then, the 2011, 2014, and 2016 World Series have gone the full seven games.
- A₄ The winner of the World Series championship is determined through a best-of-seven playoff, and the winning team is awarded the Commissioner's Trophy.
- A5 The Houston Astros won the 2017 World Series in 7 games against the Los Angeles Dodgers on November 1st, 2017, winning their first World Series since their creation in 1962.

RoBERTa

- **Q** where are trigger points located in the body
- A1 Myofascial pain is associated with muscle tenderness that arises from trigger points, focal points of tenderness, a few millimeters in diameter, found at multiple sites in a muscle and the fascia of muscle tissue.
- A2 Myofascial trigger points, also known as trigger points, are described as hyperirritable spots in the fascia surrounding skeletal muscle.
- **A**₃ Trigger points form only in muscles.
- A₄ These in turn can pull on tendons and ligaments associated with the muscle and can cause pain deep within a joint where there are no muscles.
- A5 They form as a local contraction in a small number of muscle fibers in a larger muscle or muscle bundle.

Table 4: Some qualitative examples from ASNQ test set where our ELECTRA and RoBERTa models with DSLC contextual continuous pre-training were able to rank the correct candidate in the top position while the contextual baselines failed. The answer candidates are shown ranked by the ordering produced by the contextual baselines. Other positive candidates answers are colored in light green.

I Discussion of Limitations

Our proposed pre-training approaches require ac-846 cess to large GPU resources (pre-training is per-847 formed on 350M training samples for large lan-848 guage models containing 100's of millions of pa-849 rameters). Additionally, the pre-training takes a 850 long time duration to finish (2-3 days even on 8 851 NVIDIA A100 GPUs), which highlights that this 852 procedure cannot easily be re-done with newer data 853 being made available in an online setting. How-854 ever the benefit of our approach is that once the 855 pre-training is complete, our released model check-856 points can be directly fine-tuned (even on smaller 857 target datasets) for the downstream contextual AS2 858 task. For the experiments in this paper, we only con-859 sider datasets from the English language, however 860 we conjecture that our techniques should work sim-861 ilarly for languages with limited morphology, like 862 English. Finally, we believe that the three proposed 863 objectives could be better combined in a multi-task 864 training scenario where the model has to jointly 865 predict the task and the label. At the moment, we 866 leave this as a future research direction.

845