

ViQA-COVID: COVID-19 Machine Reading Comprehension Dataset for Vietnamese

Anonymous ACL submission

Abstract

After two years of appearance, COVID-19 has negatively affected people and normal life around the world. As in November 2021, there are more than 250 million cases and five million deaths worldwide (including nearly one million cases and over twenty-two thousand deaths in Vietnam). Economy and society are both severely affected. The variant of COVID-19, Delta, has broken disease prevention measures of countries and rapidly increased number of infections. Resources overloading in treatment and epidemics prevention is happening all over the world. It can be seen that, application of artificial intelligence (AI) to support people at this time is extremely necessary. There have been many studies applying AI to prevent COVID-19 which are extremely useful, and studies on machine reading comprehension (MRC) are also in it. Realizing that, we created the first MRC dataset about COVID-19 for Vietnamese: ViQA-COVID and can be used to build models and systems, contributing to disease prevention. Besides, ViQA-COVID is also the first multi-span extraction MRC dataset for Vietnamese, we hope that it can contribute to promoting MRC studies in Vietnamese and multilingual.

1 Introduction

Delta - a so-far-most-dangerous variant of SARS-CoV-2 has shown its danger in recent months. Specifically, on average, each day there are around five hundreds thousands new cases and around ten thousands deaths worldwide. The uncontrollably rapid spread leads to the overwhelming of resources in disease prevention: medical staff, medical equipment manufacturing workers, data analysts, anti-epidemic support teams, etc. In the long run, this will have serious economic, social, as well as human impacts.

The application of intelligent and automated systems such as artificial intelligence and machine

learning models to assist and replace work for humans is essential: to help reduce the work load on staffs and hence, can help the fight against the pandemic more effective. Many researches on artificial intelligence have been applied and brought positive effects, greatly benefiting society. Among them, natural language processing (NLP) systems have worked extremely well and there will be many more applications that can be deployed in the future. Specifically, NLP applications in COVID-19 prevention can be mentioned as: information extraction; patient and location information query; automatic documents summarization; automatic question and answer about the pandemic system; epidemic related sentiment analysis, etc.

In the above application systems, MRC systems are used in many applications, but most of them only work in high-resources languages such as English, Chinese. It can be seen that dataset plays a very important role in the development of NLP systems in particular and machine learning in general. Thus, low-resource languages enrichment is the first step to create a MRC system for these languages. With that in mind, we introduce ViQA-COVID, a multi-span extraction MRC dataset related to COVID-19 for Vietnamese with the desire to develop Vietnamese and multilingual MRC systems and help in fighting against COVID-19. Dataset comprises question-answer pairs based on CDC case reports, assigned to COVID-19 prevention teams and news on Vietnam's reputation online newspapers.

In the next section, we review on the related works. Section 3 presents about datasets, statistics and annotation process. Section 4 is devoted for experiments set up. The results and benchmark are described in Section 5. Section 6 summarizes the study and presents further research directions.

2 Related Work

In recent years, COVID-19 has spurred research in many fields especially in AI related ones. In the field of computer vision, researchers (Wang et al., 2020a) designed COVID-Net to detect COVID-19 cases from chest X-ray (CXR) images and introduced COVIDx, a dataset consisting of 13,975 CXR images across 13,870 patient cases. In (Wang et al., 2020b), three masked face datasets: Masked Face Detection Dataset (MFDD), Real-world Masked Face Recognition Dataset (RMFRD), and Simulated Masked Face Recognition Dataset (SMFRD) that helped a lot in detecting and reminding people to wear masks (one of the most effective measures to prevent covid-19), are introduced. The image editing approach and datasets: Correctly Masked Face Dataset (CMFD), Incorrectly Masked Face Dataset (IMFD), as well as their combination - masked face detection (MaskedFace-Net) are introduced in (Cabani et al., 2020). MaskedFace-Net has been applied to detect whether people are wearing masks and wearing them correctly.

In the field of NLP, COVID-QA (Möller et al., 2020) is an MRC dataset consisting of 2,019 pairs of questions - answers labeled by experts, with data sources collected from COVID-19. COVID-QA is widely used in evaluating MRC tasks and applied to tasks related to COVID-19. CovidQA (Tang et al., 2020) is one of the first Question Answering datasets, consisting of pairs of questions - articles and answers that are articles related to the question. CovidDialog (Ju et al., 2020) provides a dataset including doctor-patient conversations (603 consultations and 1,232 utterances in English and 399 consultations and 8,440 utterances in Chinese). Using CovidDialog, researchers (Zeng et al., 2020) have developed a medical dialogue system to provide information related to the pandemic. Phoner_COVID (Truong et al., 2021) is a Vietnamese NER dataset about COVID-19 which defined 10 entities related to COVID-19 patients information. In addition, there are many research works that have been highly applicable and have greatly supported countries in preventing COVID-19.

With the rapid development of NLP in Vietnam, many new datasets have been introduced. From collecting 174 articles on the Vietnamese Wiki and through a five-phase annotate process, UIT-ViQuAD (Nguyen et al., 2020a) was created

with more than 23,000 question-answer pairs based on 5,109 passages. UIT-ViQuAD is a single span-extraction MRC datasets widely used in span-extraction MRC task Vietnamese besides UIT-ViNewsQA (Nguyen et al., 2020b), a dataset in healthcare domain consisting of 22,057 question-answer pairs based on 4,416 articles health report. In addition, ViMMRC (Nguyen et al., 2020c) is a multiple-choice dataset and includes 2,783 multiple-choice questions based on 417 Vietnamese texts. With the task of sentence extraction-based MRC, UIT-ViWikiQA (Do et al., 2021) is the first Vietnamese sentence extraction-based MRC dataset, created from converting the UIT-ViQuAD dataset. UIT-ViWikiQA includes 23,074 question-answer pairs, based on 5,109 passages.

In addition to the studies on COVID-19 and MRC datasets for Vietnamese, we also consulted other famous MRC datasets such as: SQuAD1.1 (Rajpurkar et al., 2016), SQuAD2.0 (Rajpurkar et al., 2018), GLUE (Wang et al., 2018), SuperGLUE (Wang et al., 2019), MASH-QA (Zhu et al., 2020), QUOREF (Dasigi et al., 2019) and DROP (Dua et al., 2019).

The above studies helped us to complete our research.

3 Dataset

In this section, we will describe ViQA-COVID in detail, annotation processing and statistics about the dataset.

Being supplied data on cases, reflections and frequently asked people's questions, from CDC Vietnam, we support them in data processing to make statistical reports, forecast epidemics, extract patient's information and collect people's feedback. This work is essential in helping the understand the situation and make decisions to prevent COVID-19. However, it is extremely complicated to process data on some hundreds thousands of cases as well as complaints, questions from people every day by hand. Unfixed-form data and complexity of Vietnamese make it difficult to handle with rule-based approach. With the development of deep learning, it can solve the above problems. Based on deep learning, a MRC system can correctly return information based on passages and questions. For example: From the patient's epidemiological information, the medical team asks: "*Who has the COVID-19 patient been in contact with?*". MRC system can answer correctly and the medical team can isolate and treat

<p>Passage: Vietnamese: Ngày 12/8, 13/8, 14/8/2020, bệnh nhân chỉ ở trong phòng thuộc khu cách ly tập trung dành cho nhân viên y tế. Ngày 15/8/2020, bệnh nhân được lấy mẫu xét nghiệm dịch hầu họng (lần 3), ngày 16/8/2020 bệnh nhân có kết quả (+) với vi rút SARS-CoV-2. Hiện tại, bệnh nhân được cách ly và điều trị tại Bệnh viện Phổi Đà Nẵng. English: On 12/8, 13/8, 14/8/2020, patient was in a room in the concentrated isolation area for medical staff. On 15/8/2020, patient was sampled for oropharyngeal fluid testing (3rd time), on 16/8/2020, patient had (+) result with SARS-CoV-2 virus. Currently, patient is isolated and treated at Da Nang Lung Hospital. Question: Bệnh nhân có kết quả dương tính vào ngày nào? (What date did the patient test positive?) Answer: 16/8/2020</p>
<p>Passage: Vietnamese: Thông tin dịch tễ: khoảng 07 giờ 00 ngày 26/7/2020, bệnh nhân trở về nhà và tiếp xúc với những người trong gia đình. Khoảng 07 giờ 00 ngày 27/7/2020, bệnh nhân được cách ly tại Bệnh viện đến ngày 02/8/2020. Sáng ngày 03/8/2020, tại Bệnh viện Đà Nẵng, bệnh nhân được lấy mẫu xét nghiệm dịch hầu họng (lần 2) và có kết quả (+) với vi rút SARS-CoV-2. Bệnh nhân ở cùng phòng với anh Đ.T (bảo vệ Bệnh viện Đà Nẵng). English: Epidemiological information: around 7:00 am on 26/7/2020, patient returned home and contacted with family members. Around 7:00 am on 27/7/2020, patient was isolated at Hospital until 02/8/2020. On the morning 03/8/2020, at Da Nang Hospital, patient was sampled oropharyngeal fluid testing (2nd time) and got a (+) result for SARS-CoV-2 virus. Patient was in the same room with D.T (security guard Da Nang Hospital). Question: Bệnh nhân đã tiếp xúc với những ai? (Who has the patient been in contact with?) Answer: những người trong gia đình, anh Đ.T (bảo vệ Bệnh viện Đà Nẵng) (family members, D.T (security guard Da Nang Hospital))</p>

Figure 1: Examples include passage, question and answer from ViQA-COVID. Bold words in passage are answers

those people quickly. In addition, MRC system can help answer people’s questions about disease, policies, ways to prevent COVID-19, and so on. To be able to achieve those applications, MRC system needs to train with MRC datasets. Therefore, we created ViQA-COVID as training data for such system. Figure 1 shows examples from ViQA-COVID, including single-span and multi-span questions.

3.1 Annotation

Annotation includes following phases:

- With limited time and resource, annotating all the data is not possible. Therefore, we reviewed and chose report cases that are as informative and structurally diverse as possible. Data was encrypted sensitive information, corrected typing and grammar errors. After data cleaning, we collected a total of 537 passages.
- Data is manually annotated. Question-answer pairs in ViQA-COVID are based on the information CDC needs to support patients and prevent diseases, as well as questions from people about the epidemic situation. For example: "What places have patients been to?", "Where are the epidemic locations that I need to be aware of?", etc. Annotators will read each passage, create questions and mark spans for corresponding answers (a answer can include multi-span). Questions are diversified

Question Types	Question Words
What (19.3%)	là gì (10.5%)
Where (17.2%)	đâu (7.3%)
When (36.6%)	ngày nào (10.4%)
Who (8.4%)	ai (6.2%)
How (3.3%)	thế nào (1.1%)
How many (10.2%)	bao nhiêu (9.6%)

Figure 2: Question types and questions words distribution in ViQA-COVID

- and avoiding duplication. Question-answer pairs are cross-checked to eliminate errors. 209 210
- We have proceeded to collect more data from reputable online portals and online news sites to diversify dataset. This data is also reviewed, manually annotated, and cross-checked. 211 212 213 214
- We reviewed and cross-checked again to complete ViQA-COVID dataset. 215 216

3.2 Statistics

ViQA-COVID after completion has a total of 6,444 question-answer pairs based on 537 passages. To our knowledge, ViQA-COVID is the first multi-span extraction MRC dataset on COVID for Vietnamese. Details of the statistics are in Table 1. It can be seen that, because ViQA-COVID is a

	Train	Dev.	Test
Number of passages	284	139	114
Number of questions	3408	1668	1368
Average passage length	336.8	269.1	252.7
Average question length	11.2	9.5	11.1
Passage vocabulary size	6659	3882	3089
Question vocabulary size	1071	606	601
Number of multi-span answers (%)	712 (20.9)	351 (21.0)	291 (21.3)
Number of single-span answers (%)	2288 (67.1)	1147 (68.8)	927 (67.8)
Number of non-span answer (%)	408 (12.0)	170 (10.2)	150 (10.9)

Table 1: ViQA-COVID overview

domain-specific dataset (COVID-19 and Health), the vocab size is not too large. In addition, the percentage of multi-span answers is quite high compared to most multi-span MRC datasets, around 20%.

We statistic question types in the dataset as follows: What: 19.3%, How: 3.33%, How many: 10.2%, Where: 17.16%, When: 36.61%, Who: 8.38%, Others: 5.02%. Like many others languages, there may be some variations in Vietnamese each type of question. Statistical description of question words in ViQA-COVID is shown in Figure 2.

4 Experiments

In this section, we present experiments with the state-of-the-art MRC models on ViQA-COVID.

4.1 Models

Since BERT (Devlin et al., 2019) - a pretrained model using Transformer (Vaswani et al., 2017) architecture appeared in 2019, it has created a strong development in the field of natural language processing. State-of-the-art performance on NLP tasks increased rapidly thanks to improved models from both BERT and the Transformer architecture. It can be said that they are the two main factors that create a new era for NLP. In this experimental part, we used variants of BERT to evaluate on ViQA-COVID. These models have achieved state-of-the-art results on many MRC tasks.

- **mBERT**: twelve layers with twelve self-attention heads BERT is trained on multi-lingual datasets (including Vietnamese). Since its launch in 2019, mBERT has performed very well in multi-lingual MRC and NLP tasks.

Passage Length	Train	Dev.	Test	Total
< 128 tokens	0	0	2	2
128 - 256 tokens	12	1	2	15
256 - 384 tokens	25	11	8	44
384 - 512 tokens	38	20	18	76
≥ 512 tokens	260	119	96	475

Table 2: Passage length and Question length statistics

- **XLM-R** (Conneau et al., 2020): based on RoBERTa (Liu et al., 2019) - an optimal BERT-based approach, XLM-R was trained on over two terabytes of cleaned Common-Crawl (Wenzek et al., 2019) data in 100 languages. XLM-R outperformed mBERT in many cross-lingual benchmarks and other tasks. We evaluated two model - XLM-R_{base}: 12 layers with 8 self-attention heads and XLM-R_{large}: 24 layers with 16 self-attention heads.
- **PhoBERT** (Nguyen and Nguyen, 2020): based on RoBERTa, PhoBERT is a Vietnamese model which improved the state-of-the-art many Vietnamese NLP tasks. PhoBERT is trained on over 20 gigabytes of word-level data (while other models train with syllable data). We also evaluated two models: PhoBERT_{base}: 12 layers with 12 self-attention heads and PhoBERT_{large}: 24 layers with 16 self-attention heads

4.2 Input Processing

Statistics from Table 2 show that most passages are in excess of 512 tokens in length. Whereas maximum length of models' input feature is 512 tokens. To deal with very long passage, we split one example into input features, each of the length

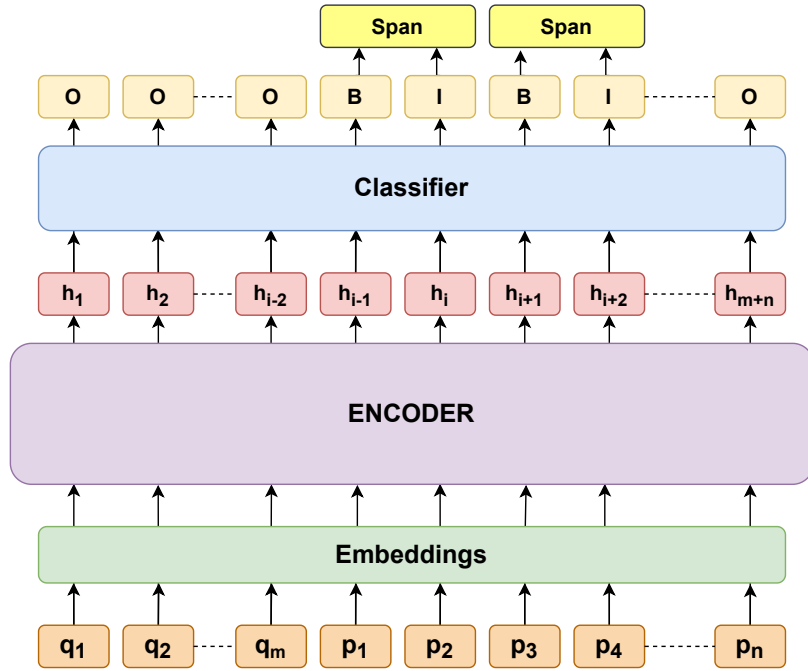


Figure 3: Illustrating the sequence tagging approach for multi-span questions. In which, $\{q_j\}_{j=1}^m$ are question tokens, $\{p_k\}_{k=1}^n$ are passage tokens and $\{h_i\}_{i=1}^{m+n}$ are the contextualized representations of the input tokens.

is shorter than model’s maximum length. In case the answer lies at the position that long passage was split, we create an overlap feature between two features (controlled by stride parameter).

PhoBERT is trained with both syllable-level and word-level tokens. Unlike English, words in Vietnamese can be compound words, i.e. one word with single meaning may be a combinations of two or more single words and in most of the cases, the meaning of the compound word is very different from their components. Thus, input sentences are segmented by word segmentation which can represent them in either syllable or word-level. Therefore, word segmentation joins syllables with a "_" sign to indicate it’s a word and makes sentences have clearer meanings. With that idea, PhoBERT outperformed XLM-R in many Vietnamese-specific NLP tasks. In our experiment, we use RDRSegmenter (Nguyen et al., 2018) from VnCoreNLP (Vu et al., 2018) as word and sentence segmentation.

4.3 Multi-span Approach

For the BERT-style models, we use sequence tagging approach (Segal et al., 2020) for multi-span questions. Instead of predicting start and end prob-

abilities like single-span questions, we predict the tag for each token. The familiar tags used are B, I, O, where B is the starting token and I is the subsequent token in output span, O is the token that is not part of an output span. Multi-span can be extracted based on B, O tokens. Figure 3 illustrates this approach in detail.

4.4 Training

BERT-style models have maximum input features length of 384 (PhoBERT of 256) with stride parameter of 128. We fine-tuned models with AdamW (Loshchilov and Hutter, 2019), weight decay of 0.01, learning rate of 5e-5 and batch size of 32, in 30 training epochs on a NVIDIA Tesla P100 GPU via Google Colaboratory. Task performance was evaluated after each epoch on the development set.

5 Results

We evaluated models’ performance on ViQA-COVID using exact match (EM) and F1-score. Results are shown in Table 3. In which, XLM-R_{large} outperforms other models with 83.37% F1-score and 68.82% EM on development set and 85.97% F1-score and 72.00% EM on test set. We also evaluated the performance of the models on single-span

Model	Dev.						Test					
	Single-Span		Multi-Span		All		Single-Span		Multi-Span		All	
	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1
mBERT	45.10	51.09	30.52	61.81	40.83	54.44	46.28	51.14	37.64	65.98	43.49	55.96
PhoBERT _{base}	61.86	72.73	30.59	54.12	51.37	66.49	54.90	74.39	34.99	60.87	54.89	70.01
PhoBERT _{large}	62.13	72.48	32.74	56.70	52.28	67.19	64.65	74.21	37.25	62.14	55.77	70.30
XLM-R _{base}	78.90	83.32	33.20	71.95	64.62	79.77	81.23	85.13	41.27	77.83	68.34	82.78
XLM-R _{large}	82.74	86.79	38.20	75.84	68.82	83.37	85.11	89.24	44.44	79.10	72.00	85.97

Table 3: Performances on development set and test set

Question Types	Dev. errors	Test errors
When	173	163
Where	98	84
Who	83	72
Others	135	95

Table 4: Question type errors on development set and test set

and multi-span answers. The models are quite accurate in predicting single-span answers but still have difficulties with multi-span answers, especially in terms of exact matching. Overall, XLM-R_{large} performed quite well and the difficulty of ViQA-COVID is not too hard when compare to other MRC datasets.

5.1 Error analyst

Through empirical analysis with the best model XLM-R_{large}, we have counted the number of incorrect answers in the development set and test set. The development set has 489/1,668 incorrect answers of which 162 multi-span, 246 single-span and 81 non-span answers. The test set has 414/1,368 incorrect answers of which 141 multi-span, 203 single-span, and 70 non-span answers. We divide these errors into four groups:

- The first group consists of answers that have the correct number of spans but have an excess or lack of words. The cases are mostly long addresses or time periods (e.g. “20/5/2020 to 30/5/2020” but the model can only predict “20/5” or “30/5”). These are also common mistakes in sequence tagging models.
- The second group includes answers that have an excess or lack of span. Mainly occurs when encountering questions about many places or about many people. For example: answering a question that lists people who have been in contact with the patient but also lists those who have not.

- The third group are completely incorrect answers (answers that have no correct span), often occurring in passages having a lot of noise. For example: Patient’s epidemiological report contains multiple dates, including dates of admission. When answering the question about the date of admission for COVID-19 infection, the model easily mistakenly answered to the date the patient was hospitalized for another illness because of the same keyword “admission”.
- The fourth group includes incorrect answers on other types of questions.

The statistics of the incorrect answers are shown in Table 4.

6 Conclusion

In this study, we introduced ViQA-COVID, the first multi-span MRC dataset about COVID-19 for Vietnamese. Our dataset consists of 6,444 question-answer pairs based on 537 passages related to COVID-19. We also experimented with different the state-of-the-art MRC models on ViQA-COVID. The results show that, XLM-R_{large} outperforms other models with 83.37% F1-score and 68.82% EM on development set and 85.97% F1-score and 72.00% EM on test set. We hope that our dataset will contribute to the prevention of COVID-19 as well as the development of NLP for Vietnamese and multilingual.

References

- Adnane Cabani, Karim Hammoudi, Halim Benhabiles, and Mahmoud Melkemi. 2020. [Maskedface-net – a dataset of correctly/incorrectly masked face images in the context of covid-19](#). *Smart Health*.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised](#)

404	cross-lingual representation learning at scale. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8440–8451, Online. Association for Computational Linguistics.	459
405		460
406		461
407		462
408		
409	Pradeep Dasigi, Nelson F. Liu, Ana Marasović, Noah A. Smith, and Matt Gardner. 2019. Quoref: A reading comprehension dataset with questions requiring coreferential reasoning . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 5925–5932, Hong Kong, China. Association for Computational Linguistics.	463
410		464
411		465
412		466
413		467
414		468
415		
416		469
417		470
		471
		472
418	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding . In <i>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)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	473
419		474
420		475
421		
422		476
423		477
424		478
425		479
426		
427	Phong Nguyen-Thuan Do, Nhat Duy Nguyen, Tin Van Huynh, Kiet Van Nguyen, Anh Gia-Tuan Nguyen, and Ngan Luu-Thuy Nguyen. 2021. Sentence extraction-based machine reading comprehension for vietnamese . <i>CoRR</i> , abs/2105.09043.	480
428		481
429		482
430		483
431		484
		485
432	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 <i>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)</i> , pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.	486
433		487
434		488
435		489
436		490
437		491
438		492
439		
440		493
441		494
		495
		496
442	Zeqian Ju, Subrato Chakravorty, Xuehai He, Shu Chen, Xingyi Yang, and Pengtao Xie. 2020. Coviddialog: Medical dialogue datasets about covid-19 . https://github.com/UCSD-AI4H/COVID-Dialogue .	497
443		498
444		
445		499
446	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 . Cite arxiv:1907.11692.	500
447		501
448		502
449		503
450		504
		505
451	Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization . In <i>International Conference on Learning Representations</i> .	506
452		507
453		508
454		509
455		
456		510
457		511
458		512
		513
		514
		515
		459
		460
		461
		462
		463
		464
		465
		466
		467
		468
		469
		470
		471
		472
		473
		474
		475
		476
		477
		478
		479
		480
		481
		482
		483
		484
		485
		486
		487
		488
		489
		490
		491
		492
		493
		494
		495
		496
		497
		498
		499
		500
		501
		502
		503
		504
		505
		506
		507
		508
		509
		510
		511
		512
		513
		514
		515

- 516 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
517 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
518 Kaiser, and Illia Polosukhin. 2017. [Attention is all
519 you need](#). In *Advances in Neural Information Pro-
520 cessing Systems*, volume 30. Curran Associates, Inc.
- 521 Thanh Vu, Dat Quoc Nguyen, Dai Quoc Nguyen, Mark
522 Dras, and Mark Johnson. 2018. [VnCoreNLP: A Viet-
523 namese natural language processing toolkit](#). In *Pro-
524 ceedings of the 2018 Conference of the North Amer-
525 ican Chapter of the Association for Computational
526 Linguistics: Demonstrations*, pages 56–60, New Or-
527 leans, Louisiana. Association for Computational Lin-
528 guistics.
- 529 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman-
530 preet Singh, Julian Michael, Felix Hill, Omer Levy,
531 and Samuel Bowman. 2019. [Superglue: A stickier
532 benchmark for general-purpose language understand-
533 ing systems](#). In *Advances in Neural Information Pro-
534 cessing Systems*, volume 32. Curran Associates, Inc.
- 535 Alex Wang, Amanpreet Singh, Julian Michael, Fe-
536elix Hill, Omer Levy, and Samuel Bowman. 2018.
537 [GLUE: A multi-task benchmark and analysis plat-
538 form for natural language understanding](#). In *Proce-
539 edings of the 2018 EMNLP Workshop BlackboxNLP:
540 Analyzing and Interpreting Neural Networks for NLP*,
541 pages 353–355, Brussels, Belgium. Association for
542 Computational Linguistics.
- 543 Linda Wang, Zhong Qiu Lin, and Alexander Wong.
544 2020a. [Covid-net: a tailored deep convolutional
545 neural network design for detection of covid-19
546 cases from chest x-ray images](#). *Scientific Reports*,
547 10(1):19549.
- 548 Zhongyuan Wang, Guangcheng Wang, Baojin Huang,
549 Zhangyang Xiong, Qi Hong, Hao Wu, Peng Yi,
550 Kui Jiang, Nanxi Wang, Yingjiao Pei, Heling Chen,
551 Yu Miao, Zhibing Huang, and Jinbi Liang. 2020b.
552 [Masked face recognition dataset and application](#).
553 *CoRR*, abs/2003.09093.
- 554 Guillaume Wenzek, Marie-Anne Lachaux, Alexis Con-
555 neau, Vishrav Chaudhary, Francisco Guzmán, Ar-
556 mand Joulin, and Edouard Grave. 2019. [Ccnet: Ex-
557 tracting high quality monolingual datasets from web
558 crawl data](#).
- 559 Guangtao Zeng, Qingyang Wu, Yichen Zhang, Zhou
560 Yu, Eric Xing, and Pengtao Xie. 2020. De-
561 velop medical dialogue systems for covid-19.
562 <https://github.com/UCSD-AI4H/COVID-Dialogue>.
- 563 Ming Zhu, Aman Ahuja, Da-Cheng Juan, Wei Wei, and
564 Chandan K. Reddy. 2020. [Question answering with
565 long multiple-span answers](#). In *Findings of the Asso-
566 ciation for Computational Linguistics: EMNLP 2020*,
567 pages 3840–3849, Online. Association for Computa-
568 tional Linguistics.